

# FloCon 2004 Proceedings

**Jul 2004**

**CERT Program**

<http://www.sei.cmu.edu>



Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE <b>JUL 2004</b>		2. REPORT TYPE		3. DATES COVERED <b>00-00-2004 to 00-00-2004</b>	
4. TITLE AND SUBTITLE <b>FloCon 2004 Proceedings</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>Carnegie Mellon University,Software Engineering Institute,Pittsburgh,PA,15213</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release; distribution unlimited</b>					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT <b>Same as Report (SAR)</b>	18. NUMBER OF PAGES <b>195</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

Copyright 2004 Carnegie Mellon University

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8721-05-C-0003 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the United States Department of Defense.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

This material has been approved for public release and unlimited distribution except as restricted below.

Internal use:\* Permission to reproduce this material and to prepare derivative works from this material for internal use is granted, provided the copyright and "No Warranty" statements are included with all reproductions and derivative works.

External use:\* This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other external and/or commercial use. Requests for permission should be directed to the Software Engineering Institute at [permission@sei.cmu.edu](mailto:permission@sei.cmu.edu).

\* These restrictions do not apply to U.S. government entities.



**Carnegie Mellon  
Software Engineering Institute**

**CERT**  
Analysis  
Center

# **Empirically Based Analysis: The DDoS Case**

**Jul 22<sup>nd</sup>, 2004**

**CERT® Analysis Center  
Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890**

*The CERT Analysis Center is part of the Software Engineering Institute. The Software Engineering Institute is sponsored by the U.S. Department of Defense.*

*© 2003 by Carnegie Mellon University*





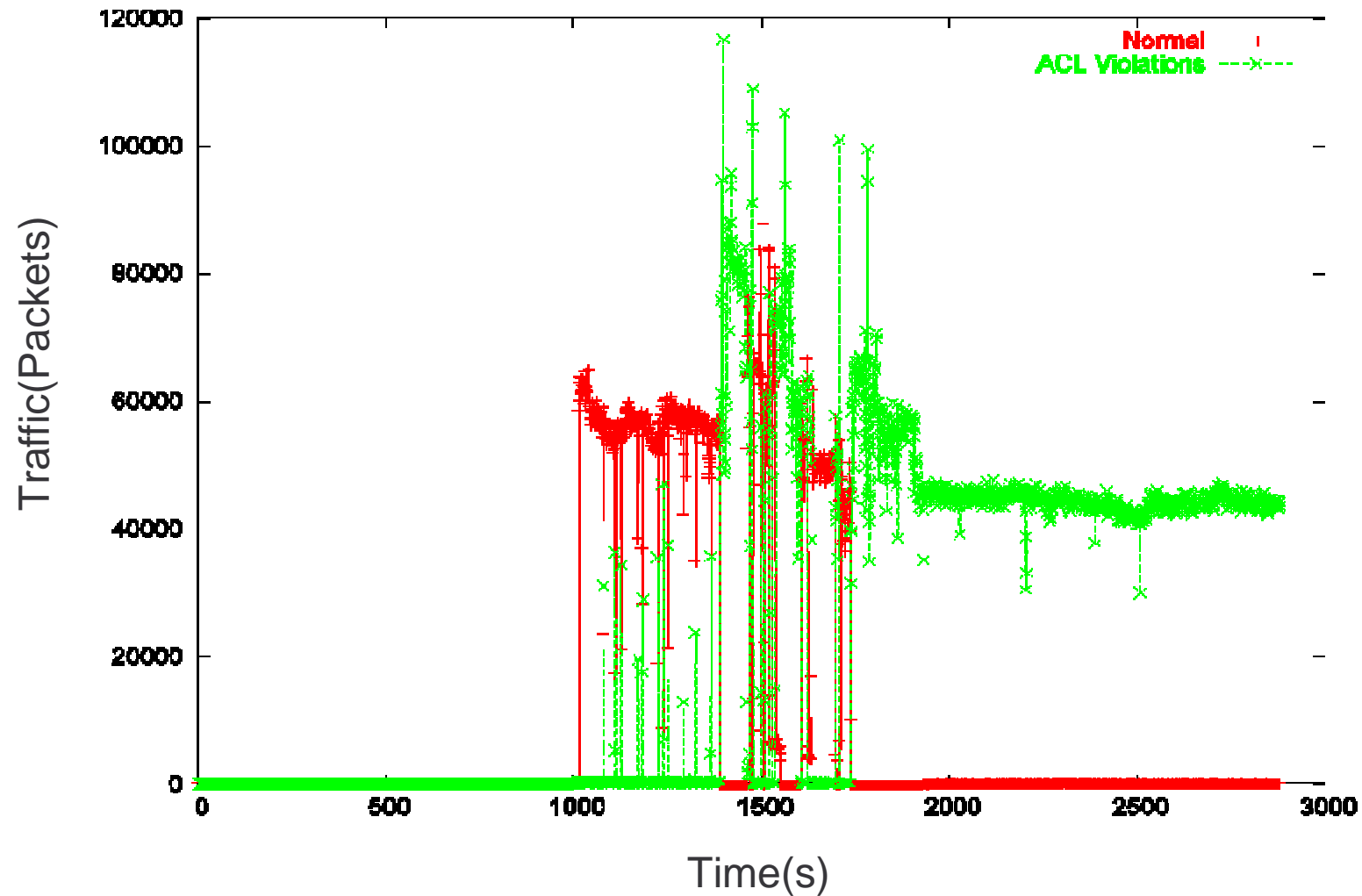


# Introduction

- Ø Access to the dataset gives us a large enough record of traffic to test hypotheses in network security.
- Ø Given this, we select and evaluate various security measures against real traffic
  - Or a reasonable facsimile thereof
- Ø One example: target resident DDoS Filters
  - Heavily constrain the problem— not considering SYN floods, smurfing, reflection attacks...



# Attacks like this





# How Do We Test?

Ø Any analysis opens a can of worms...err,  
“assumptions”

- The network constantly changes
- What is a representative host?

Ø Rerunning attacks is of debatable value

- Most of the legitimate traffic is dropped, that's what a DoS is *for*

Ø We want our results to be representative

- Test and summarize over multiple machines

Ø We want our results to be reproducible

- Depend heavily on SiLK structures and tools



# Evaluation

Ø Trained filters on 15 days of legitimate traffic

- Built a representation of IP address: volume relationship (via `rwaddrcount`)

Ø Then generated a simulated DoS

- Botnet IPs collected with `rwset`
- Normal traffic selected from another day

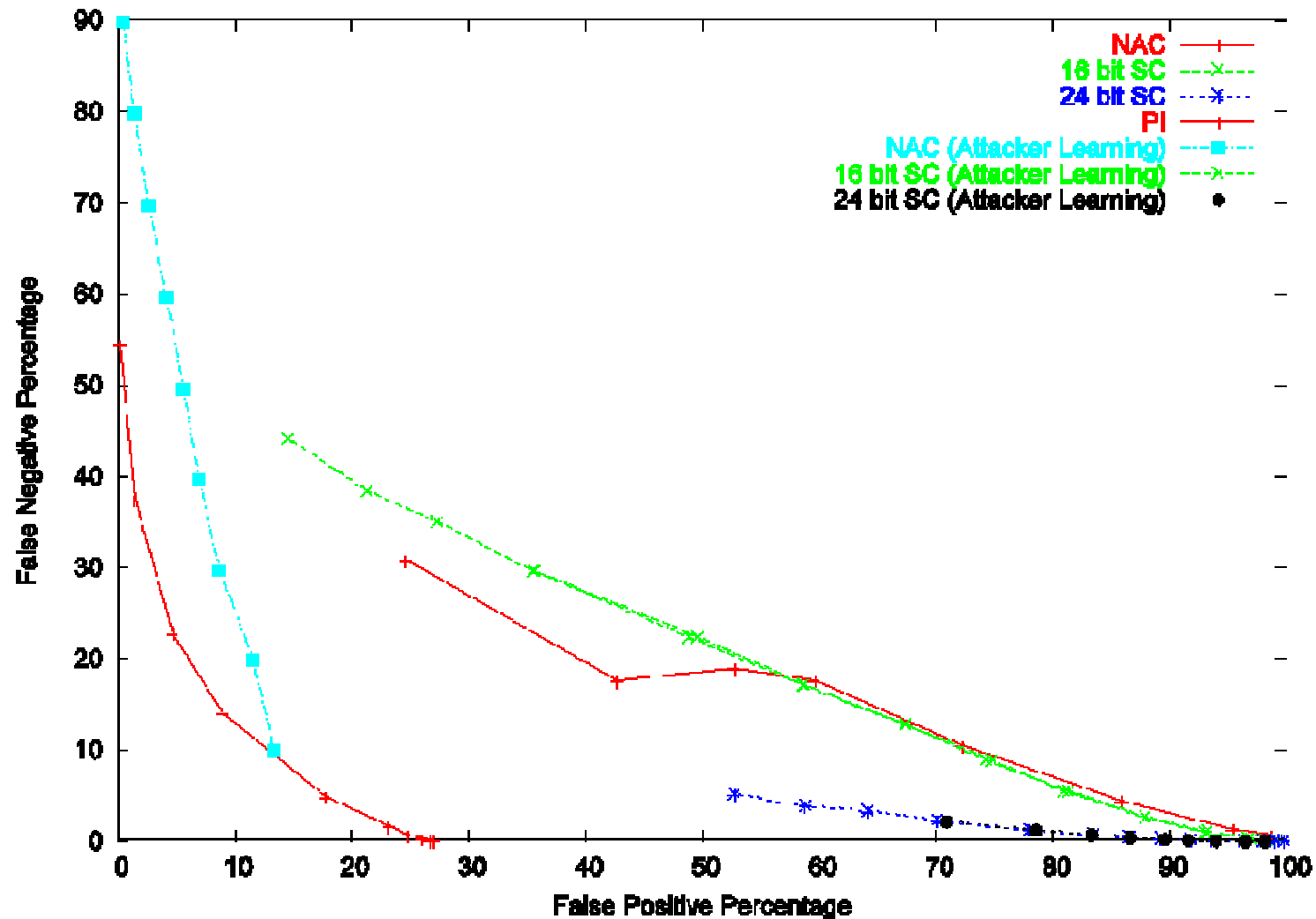
Ø Resulting traffic was then evaluated for failure rates

Ø Tested 2 types of filters:

- Clustering – groups of adjacent IP addresses
- PI – path marking approach



# DoS Filters





# Initial Observations

## Ø Two groups

- One group assumes a magic DoS Detection Oracle
  - That's the group with better results

## Ø In general, the filters don't do well

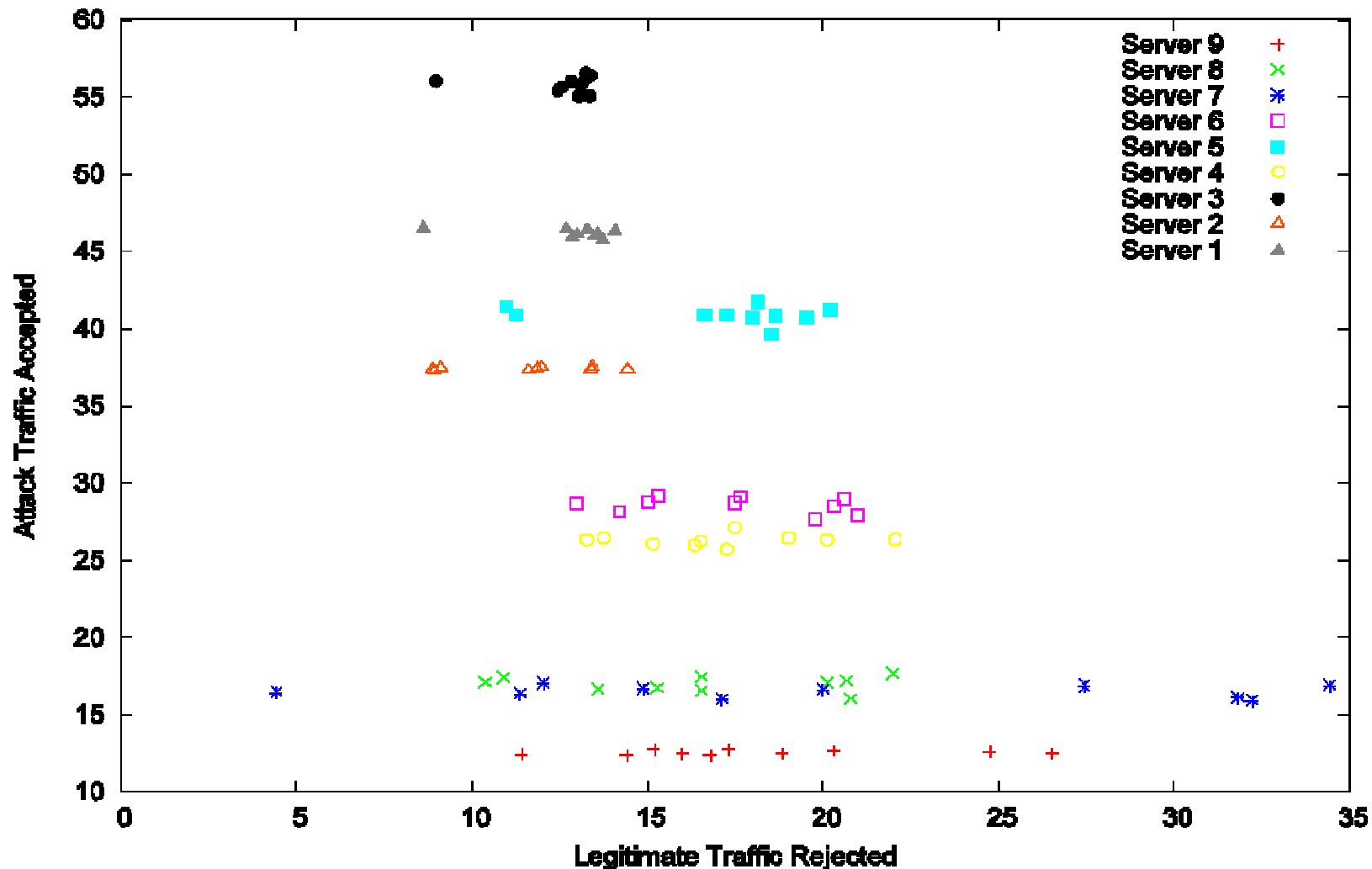
- Should we compare IP addresses, or packets?
- Is traffic different for different servers?

## Ø Let's look at one result in more depth



# One result in more depth

Comparative Failure Rates For 90% threshold, 25 Days Learning Time





# Observations

## Ø Normal traffic varies extensively

- Although it seems to vary more with “smaller” servers
- And it’s better when you look at packet counts
  - Which makes sense, given the absurd number of scanners we see.

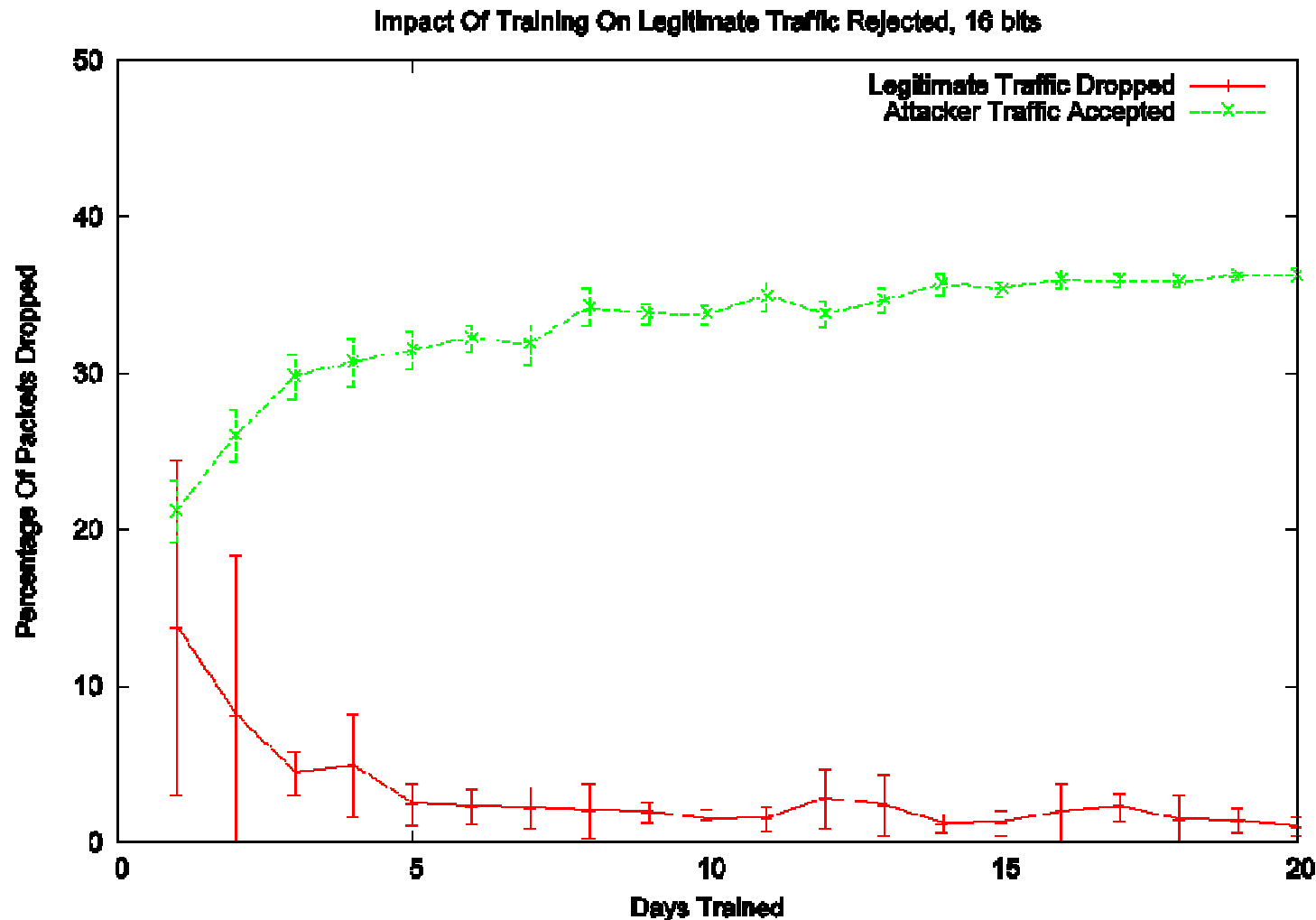
## Ø False negative rate (attackers accepted) seems to be related to server activity – the busier the higher.

- Attackers don’t vary as much





# Learning Curves – 95% threshold





# Other Observations

Ø In the majority of cases, packets are dropped because they've never been seen before

- Short learning curves – effectively no change in false positive rate after a week of learning.
- Especially true for spoofed traffic

Ø Entropy is lower than expected

- Filters that rely on spoof defense (HCF, PI) drop less than 10% of their packets because they detect a spoof



# Further Work

## Ø Exploiting our DoS attack traffic records further

- We know how the network reacts
- We know how the attack starts and ends
  - Which impacts learning curve for defenses that *only* profile the attack

## Ø Further use of other network maps

- Skitter (used for PI), &c.

## Ø Formalization of the techniques used

- Developed a matrix based approach for the final iteration
- Tools are going to be available publicly



# A Final Note

ØURL for the SiLK tools:  
<http://silkttools.sourceforge.net>

# **Statistical Methods for Flow Data**

Joseph B. Kadane  
Carnegie Mellon University  
Department of Statistics and  
Software Engineering Institute  
`kadane@stat.cmu.edu`

## **Outline**

1. The Issue
2. Bayesian Techniques
3. Advantages of Bayesian Techniques
4. Conclusion

## **1. The Issue**

- Existing logistic regression is described in Marc Kellner's presentation
- With 200 odd observations in a ten-dimensional space of explanatory variables, the data can be sparse. There are three sorts of responses to this
  - (a) reduce the space by deleting explanatory variables
  - (b) collect an order of magnitude more data
  - (c) use Bayesian methods to smooth the estimates

## 2. Bayesian Techniques

The logistic regression defines a **likelihood**, that is a probability distribution of the data (which is 1's and 0's, scans and non-scans) given the explanatory variables and the (uncertain) weights. The remaining ingredient is a **prior distribution** on the weights, found by interrogating one or more experts, in a process known as elicitation.

(For more on this see Elicitation, by Garthwaite, Kadane and O'Hagan [www.stat.cmu.edu/tr/tr808/tr808.html](http://www.stat.cmu.edu/tr/tr808/tr808.html)).

With these two ingredients, one can compute the posterior distribution on the weights. The posterior is proportional to the prior times the likelihood.



### **3. Advantages of Bayesian Techniques**

- pulls in discrepant and unreasonable estimates of weights
- posterior variances guide sophisticated statistical design of experiments in deciding what additional data to gather or elicitation to ask about
- conceptually easy to extend from binomial to multinomial data

#### **4. Conclusion**

Bayesian methods offer a reasonable way forward to make the logistic regression approach to scan data stable and operationally feasible.



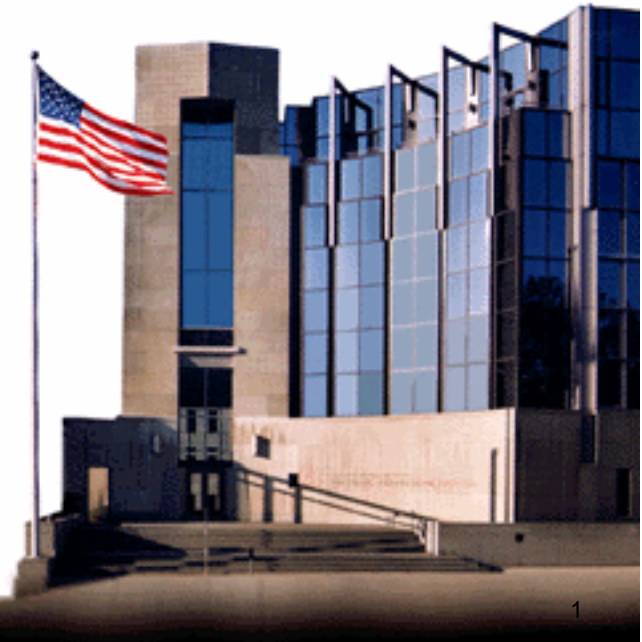
# Detection and Analysis of Scans on Very Large Networks

**FloCon 2004: Modeling Techniques Panel**  
**July 21, 2004**

**Marc Kellner**  
**Carrie Gates**

**CERT® Centers**  
**Software Engineering Institute**  
**Carnegie Mellon University**  
**Pittsburgh, PA 15213-3890**

**Sponsored by the U.S. Department of Defense**  
**© 2004 by Carnegie Mellon University**





# Needs Motivating this Approach

**A comprehensive, integrated view of scanning activity across the network(s) of interest is needed to support situational awareness.**

**A historical record of network activity is needed to detect extremely low-intensity scans.**

**A historical record of identified scans is needed to study the evolving characteristics of Internet scanning.**

**Network defenders need support to**

- **help identify higher-risk scans**
- **identify hosts at greater risk of compromise**
- **detect internal sources of scans**



# Project Components and Initial Focus

**The major thrusts of our effort are**

- **scan detection**
- **scan database**
- **analysis of scans**

**Our initial focus in scan detection is on**

- **inbound traffic**
- **single source scans**
- **using TCP protocol**

**As the scan database is populated, we will be able to commence analysis of the scans against the network(s) of interest.**



# Distinguishing Characteristics

**Unlike most scan detection approaches, ours**

- **is retrospective (not real-time)**
- **is based on flows (e.g., Cisco's NetFlow data)**
- **is multi-dimensional and extensible**
- **provides probability of traffic containing a scan**
- **supports long-term analysis of scanning activity**



# Overview of Scan Detection Steps

**Sort flow records by {source IP address (SIP), start time (stime)}**

**Identify the events (essentially clusters of traffic) for each SIP**

**Analyze each event (independently), for each SIP, to assess the probability it contains scanning activity**

**(Future) Combine traffic – not initially identified as scans – across time periods; analyze this combined traffic, for each SIP, to assess the probability it contains scanning activity**



# Scan Indicators

**We compute several scan indicators for each event. These indicators are used to compute the probability that an event contains scanning activity.**

**Indicators are computed for**

- **individual class C (/24) sub-nets (nets)**
  - **net coverage, net run length**
- **individual destination addresses (hosts)**
  - **low port coverage, low port run length**
  - **high port coverage, high port run length**
- **the event overall (event)**
  - **sub-net run length**
  - **flag combinations**
  - **packets per flow**
  - **unusual ports**





# Applying Logistic Regression

**Finally, we use the ten overall scan indicator values to determine how likely it is that an event contains scanning activity.**

**Let  $I_1, \dots, I_{10}$  represent the ten overall scan indicator values for any given event.**

**A “simple” logistic regression model using these variables predicts the probability ( $P$ ) that this event contains scanning activity as**

$$P = e^z / (1 + e^z) \quad \text{where}$$
$$z = \beta_0 + (\beta_1 * I_1) + (\beta_2 * I_2) + \dots + (\beta_{10} * I_{10})$$

**However, in order to apply this model we need to find the  $\beta$  values.**



# Model Estimation and Validation (1 of 2)

**From a dataset of 155,827 events (reflecting 56M TCP flow records) we drew two samples:**

- **for estimation of the model (120 events)**
- **for validation of the model (200 events)**

**Each of the 320 sample events were independently classified by two experts as containing scanning activity or not. This was accomplished by examining the flow records without use of any scan detection tool or the scan indicator values.**

**These classifications were used as the “gold standard” against which the model was developed and validated.**



# Model Estimation and Validation (2 of 2)

## Estimation sample:

- drawn using a stratified sampling approach
- provided to a standard logistic regression program to estimate the  $\beta$  values for the model
- results: correctly classified all 120 sample events

## Validation sample:

- drawn as a purely random sample from the dataset
- the logistic regression model was then used to classify each sample event
- results: correctly classified 197 of the 200 sample events
  - two false negatives (1.0% error rate)
  - one false positive (0.5% error rate)



# Results from the Sample Dataset

**Running the scan detection system on the full dataset of 155,827 events yielded the following:**

- **90,999 (58.40%) had prob. = 0.9 of containing scans**
- **64,288 (41.25%) had prob. = 0.1 of containing scans**
- **the remaining 540 events (0.35%) fell somewhere in between (i.e.,  $> 0.1$  and  $< 0.9$ )**

**Using the customary probability of 0.5 as the threshold for a scan, led to classifying 91,381 (58.6%) of the total events as scans.**

**As points of reference, these events**

- **contain 18,642,671 (33.1%) of the 56,344,051 TCP flow records**
- **came from 90,490 (17.7%) of the 511,602 unique source addresses sending TCP flows**



# Scan Database

**We have developed a database to record summary information about all detected scans. This will support the detection of distributed source scans and repetitive scans, as well as general analysis.**

**Vital information recorded in the database includes**

- **scan source, start time, and end time**
- **all targeted destinations (i.e., {DIP, dport} pairs)**
- **size of the scan (in number of bytes, packets, flows, unique destinations, unique dips, unique dports, /24 subnets, and duration)**
- **indicator values from the scan detection program**
- **type of scan**
- **etc.**



# Analysis of Scans

**Based on information accumulated in the scan database**

**Determine appropriate and informative scan metrics and characteristics**

**Report results such as top scan sources, top scan targets (inside protected network), top ports scanned, average scan characteristics (e.g., intensity, scan rate, duration, number of different addresses scanned, number of different ports scanned, scan flow characteristics), how frequently an average target is scanned, repetitive scan sources, etc.**



# Planned Operational Capabilities

**Phase 1: Provide capability to run scan detection and populate scan database on a routine basis**

**Phase 2: Provide basic access to scan database using a selection of pre-defined (but parameterized) queries**

**Phase 3: Provide scan information for lower tiers**

- **identify scans that included any of the addresses of interest to that lower-tier organization**
- **facilitate investigating the scan from raw flow data**

**Phase 4: Highlight hosts potentially compromised during a scan**



# Concluding Summary

**The major thrusts of our effort are**

- **scan detection**
- **scan database**
- **analysis of scans**

**Unlike most scan detection approaches, ours**

- **is retrospective (not real-time)**
- **is based on flows (e.g., Cisco's NetFlow data)**
- **is multi-dimensional and extensible**
- **provides probability of traffic containing a scan**
- **supports long-term analysis of scanning activity**

**Scan database has been designed and implemented**

**Analysis capabilities will be provided for network defenders and will support scanning research goals**



# Locality Based Analysis of Network Flows

SEI/CERT

21 July 2004

John McHugh,  
Carrie Gates, Damon Becknel

© 2004 by Carnegie Mellon  
University

## Why Locality

- Locality is an entropy based characterization that allows prediction of future behavior based on past observations.
  - It captures the degree to which the behavior of a system is regular in some sense
  - It appears to be scale free, appearing in internet, subnet, and node scale behaviors.
  - It promotes clustering allowing the use of sets and multisets to abstract group behaviors.

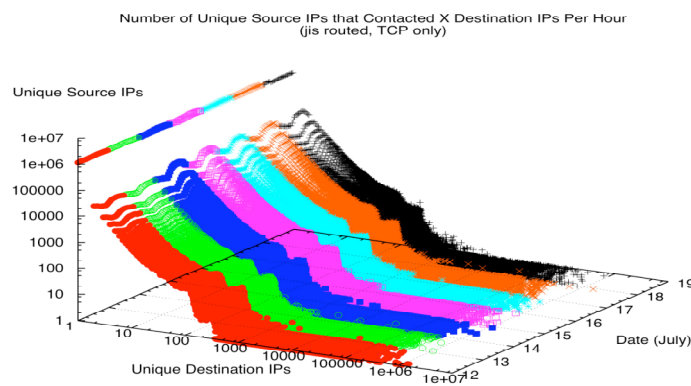
© 2004 by Carnegie Mellon  
University

## Eye Candy vs. Insight

- Locality often manifests as patterns in some space.
  - If we select the appropriate dimensions, we may achieve either understanding or puzzlement.
  - The next three pictures show persistent structure where none might be expected.
  - This can be viewed as a summary of a time series of connection matrices.
  - Graphics by Carrie Gates

© 2004 by Carnegie Mellon University

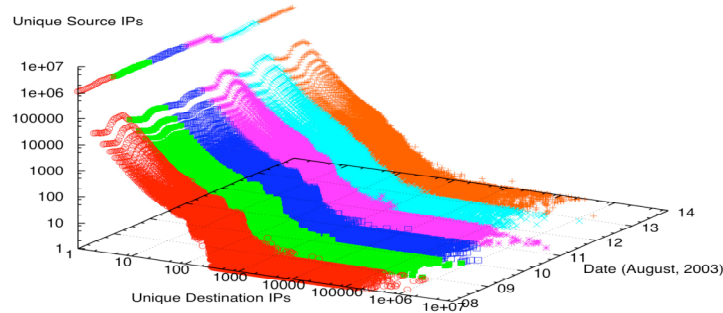
## First you see it ...



© 2004 by Carnegie Mellon University

## Then it goes away ...

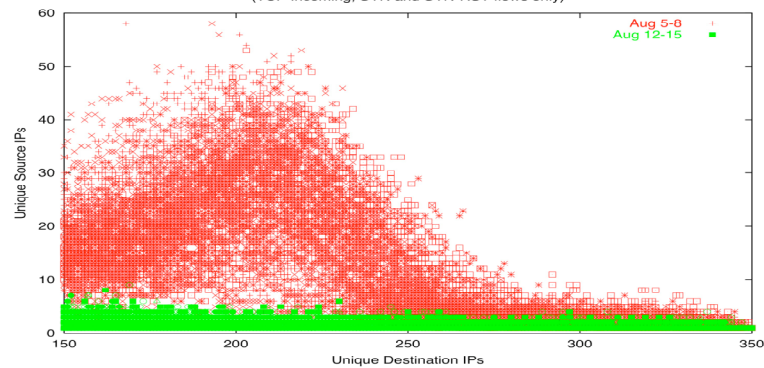
Number of Unique Source IPs that Contacted X Destination IPs Per Hour  
(jis routed, TCP only)



© 2004 by Carnegie Mellon University

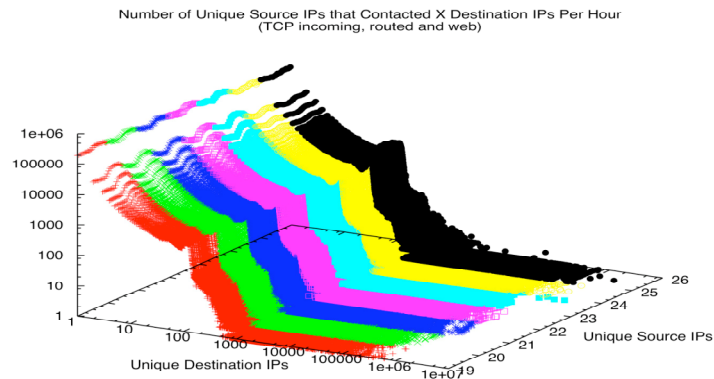
## (rather abruptly)

Number of Unique Source IPs that Contacted X Destinations Per Hour  
(TCP Incoming, SYN and SYN-RST flows only)



© 2004 by Carnegie Mellon University

## Only to return (months later).



University

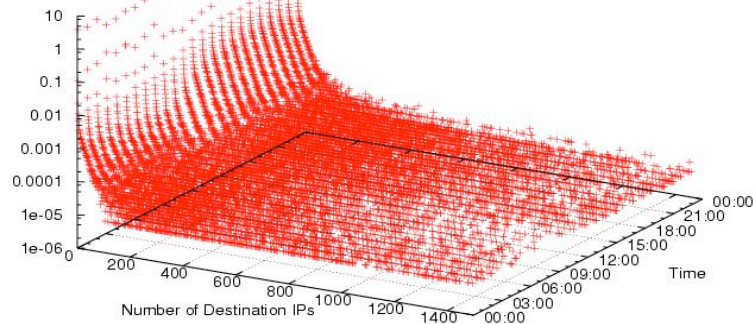
## Williamson's Locality

- Matt Williamson, late of HP Bristol, noted address locality in a 2002 ACSAC paper.
  - For browsing, last 10 IPs visited constitute an effective working set.
  - Working set violations relatively rare, bursts rarer yet.
    - Delay on violation is effective “soft” mitigator
- What is the locality of trans border data?

## Detail of Inside to Outside Day

Number of Destination IPs Contacted Per Source Over Time  
(14 January 2003, all outgoing TCP traffic, calculated on a per hour basis)

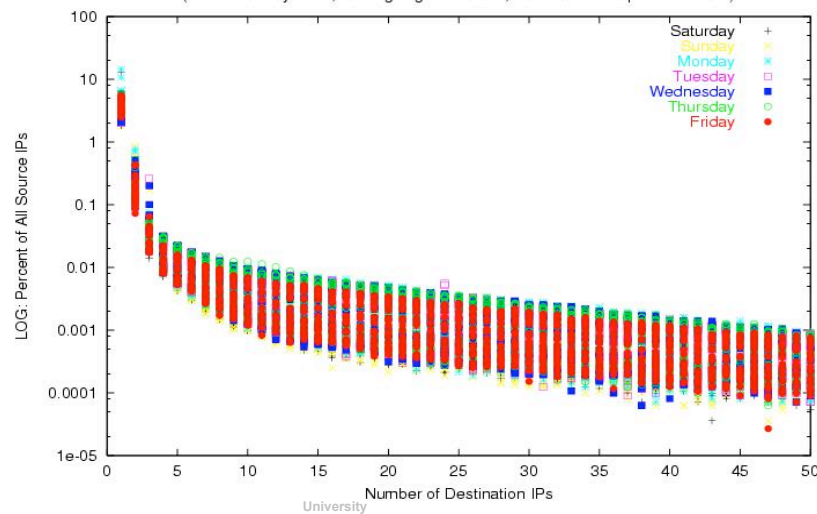
LOG: Percent of All Source IPs



University

## Weekly In/Out Locality Range

Number of Destination IPs Contacted Per Source Over Time  
(11-17 January 2003, all outgoing TCP traffic, calculated on a per hour basis)



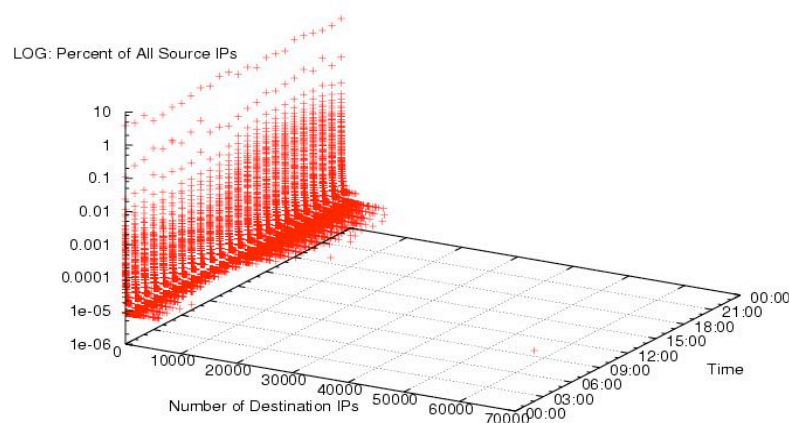
## Williamson Confirmed (mostly)

- With the caveat that we are not seeing internal connections, the vast majority of the flows arguably follow Williamson's working set model.
- As usual, there are outliers ...

© 2004 by Carnegie Mellon University

## One Day of Inside to Outside

Number of Destination IPs Contacted Per Source Over Time  
(14 January 2003, all outgoing TCP traffic, calculated on a per hour basis)



University

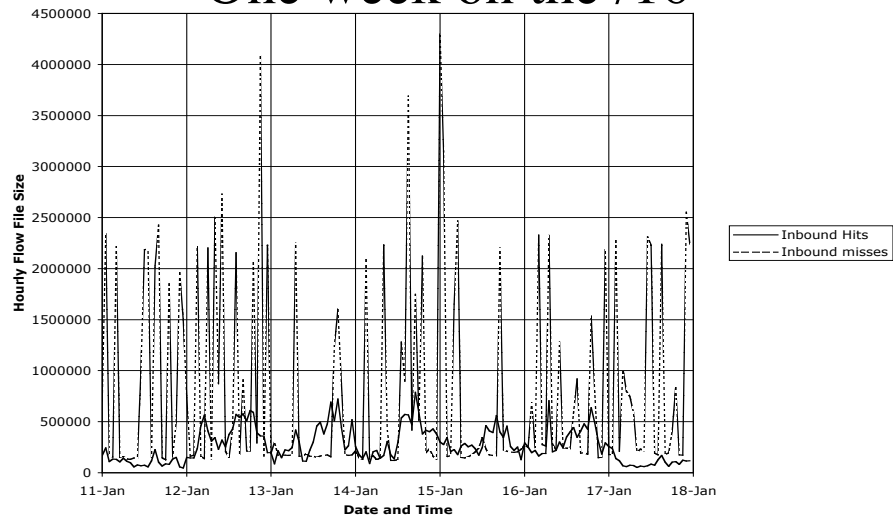
## Noise localities

- We have been characterizing modest subnets in support of the traffic generation that will be used in the DARPA DQ system evaluations.
  - Attempting to avoid mistakes of DARPA IDS evaluation.
  - Striving for a realistic noise environment, among other things.

## Crud and Noise

- In January, we observed a /16 for a week, and the whole customer net for a minute
- For the /16
  - MMM.NNN.24.x - 66 hosts    MMM.NNN.25.x - 60 hosts
  - MMM.NNN.26.x - 46 hosts    MMM.NNN.27.x - 49 hosts
  - MMM.NNN.28.x - 57 hosts    MMM.NNN.29.x - 7 hosts
  - MMM.NNN.30.x - 70 hosts    MMM.NNN.31.x - 67 hosts
  - MMM.NNN.32.x - 54 hosts    MMM.NNN.33.x - 62 hosts
  - MMM.NNN.34.x - 50 hosts    MMM.NNN.35.x - 4 hosts
  - MMM.NNN.120.x - 2 hosts    MMM.NNN.127.x - 1 host
  - MMM.NNN.140.x - 1 host    MMM.NNN.251.x - 4 hosts
  - Total 600 hosts in 16 /24s

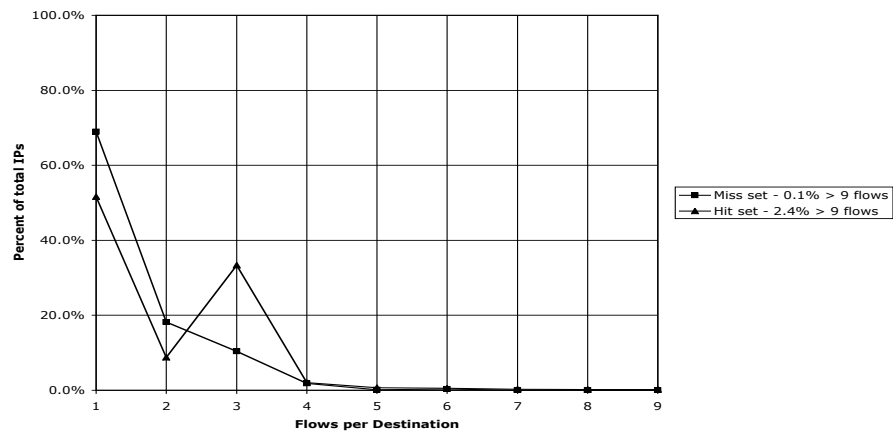
## One week on the /16



© 2004 by Carnegie Mellon University

## 1 Min sample - destinations

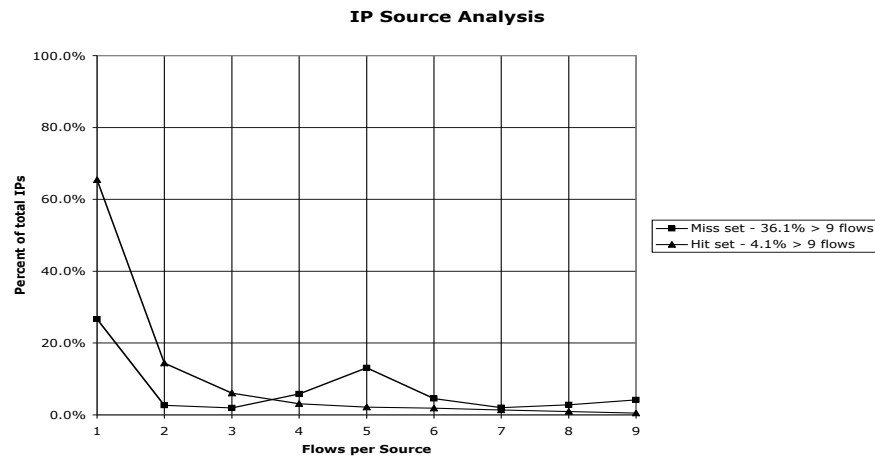
IP Destination Analysis



© 2004 by Carnegie Mellon University



# 1 Min Sample - sources



© 2004 by Carnegie Mellon University

## top 5 in 1 min sample

- Created a “bag” for source and destination addresses in the 1 minute sample. The annotated top 5 are:
- (39) `lip $ readbag --count --print jcm-tcp-s-10+.bag | sort -r -n | head`
  - 12994 AAA.BBB.068.218 - scan 4899 (Radmin)
  - 6598 CCC.DDD.209.215 - scan 7100 (X-Font)
  - 5944 EEE.FFF.125.117 - scan 20168 (Lovegate)
  - 5465 GGG.HHH.114.052 - ditto
  - 5303 III.JJJ.164.126 - scan 3127 (My doom)

© 2004 by Carnegie Mellon University

## Bottom of bag in 1 min sample

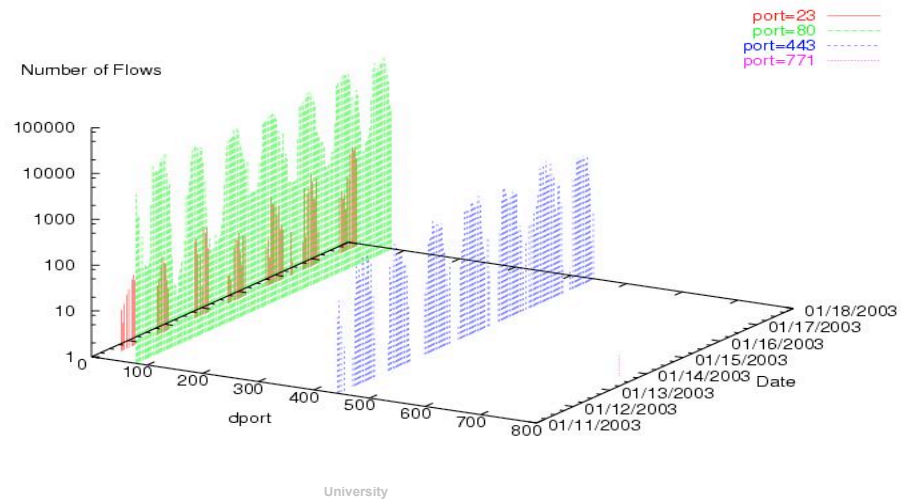
- 3335 external hosts sent exactly one TCP flow
  - SYN probes for port 8866 449 times
    - W32.Beagle.B@mm is a mass-mailing worm-back door on TCP port 8866.
  - SYN probes for port 25 are seen 271 times.
  - Most remainder are SYNs to a variety of ports, mostly with high port numbers.
  - There are a number of ACK/RST packets which are probably associated with responses to spoofed DDoS attacks.

## Individual host profiles

- These were done by Capt. Damon Becknel, USA.
  - He was looking for ways of characterizing the role of a node based on it's activity patterns
  - As usual, surprising results are sometimes observed.

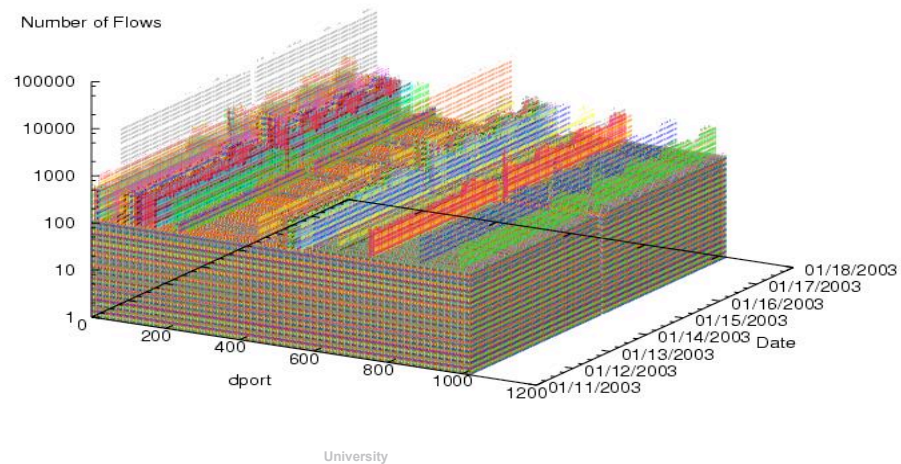
# Workstation?

Workstation? - Distribution of dport



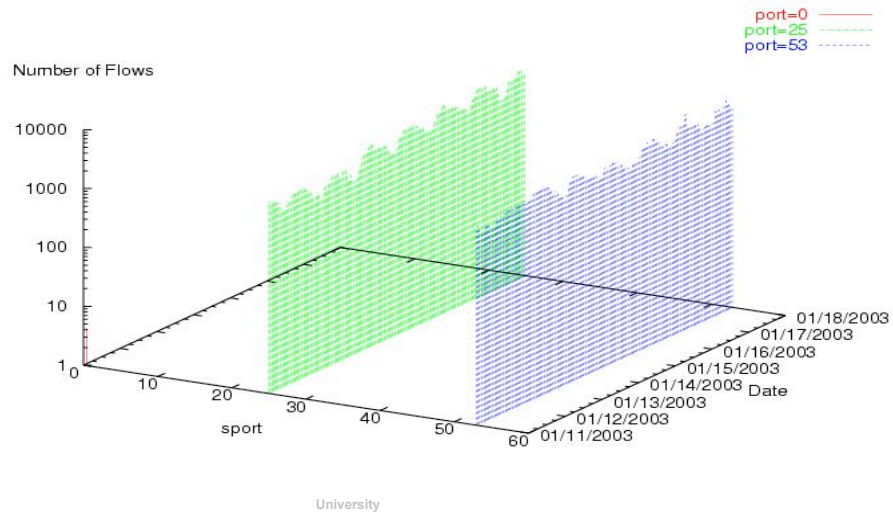
# Scanner

Scanner - Distribution of dport



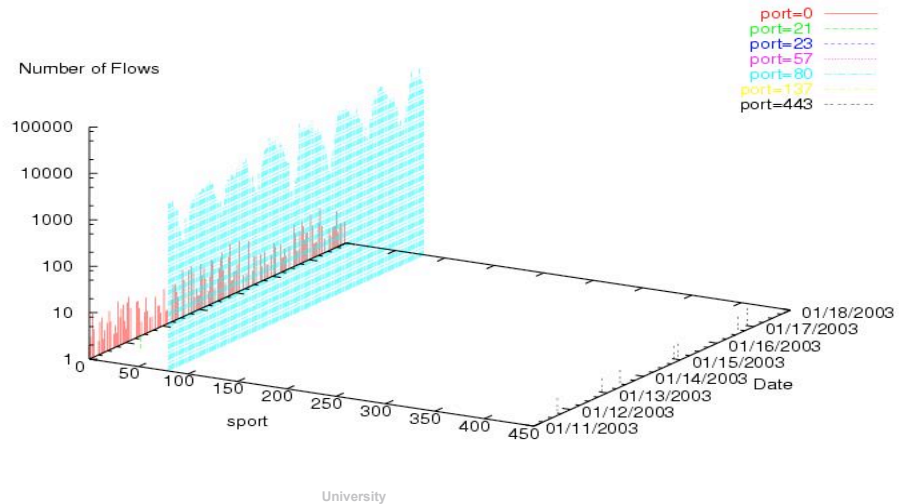
# Mail Server?

Mail Server - Distribution of sport



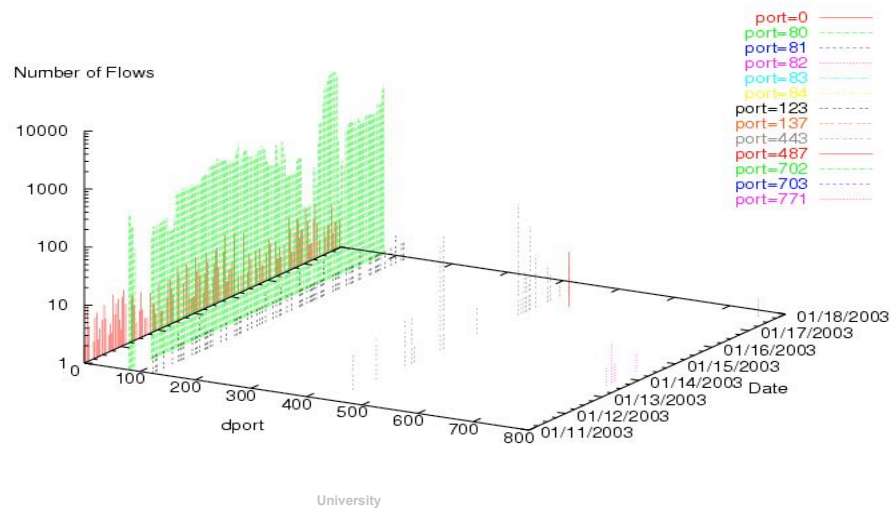
# Web Server

Web Server - Distribution of sport



# Web Server

Web Server - Distribution of dport



## Summary

- We have provided some examples of locality on a variety of scales for a variety of representations.
- It is our hope that the general notions of locality, and clustering will provide a basis for reducing the complexity of analysis.



Carnegie Mellon  
Software Engineering Institute

**CERT**  
Situational  
Awareness

# Analysis of the US-CERT DAC

Josh McNutt <[jmcnutt@cert.org](mailto:jmcnutt@cert.org)>

FloCon: Netflow Analysis Workshop

July 21, 2004

CERT® Network Situational Awareness Group  
Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

*The CERT Network Situational Awareness Group  
is part of the Software Engineering Institute.  
The Software Engineering Institute is sponsored by  
the U.S. Department of Defense.*





## Outline

---

- Data
- Graphical Displays
- Detecting Trends
- Anomaly Detection
- Roadmap



## Data

---

- **Snort**
  - Signature-based alerts
  - Pre-processor alerts
- **Origin**
  - Multiple networks of varying size
- **Volume**
  - ~30-50 million alerts per month
- **Ancillary Information**
  - Country code
  - Netblock





## IDS Data: challenges

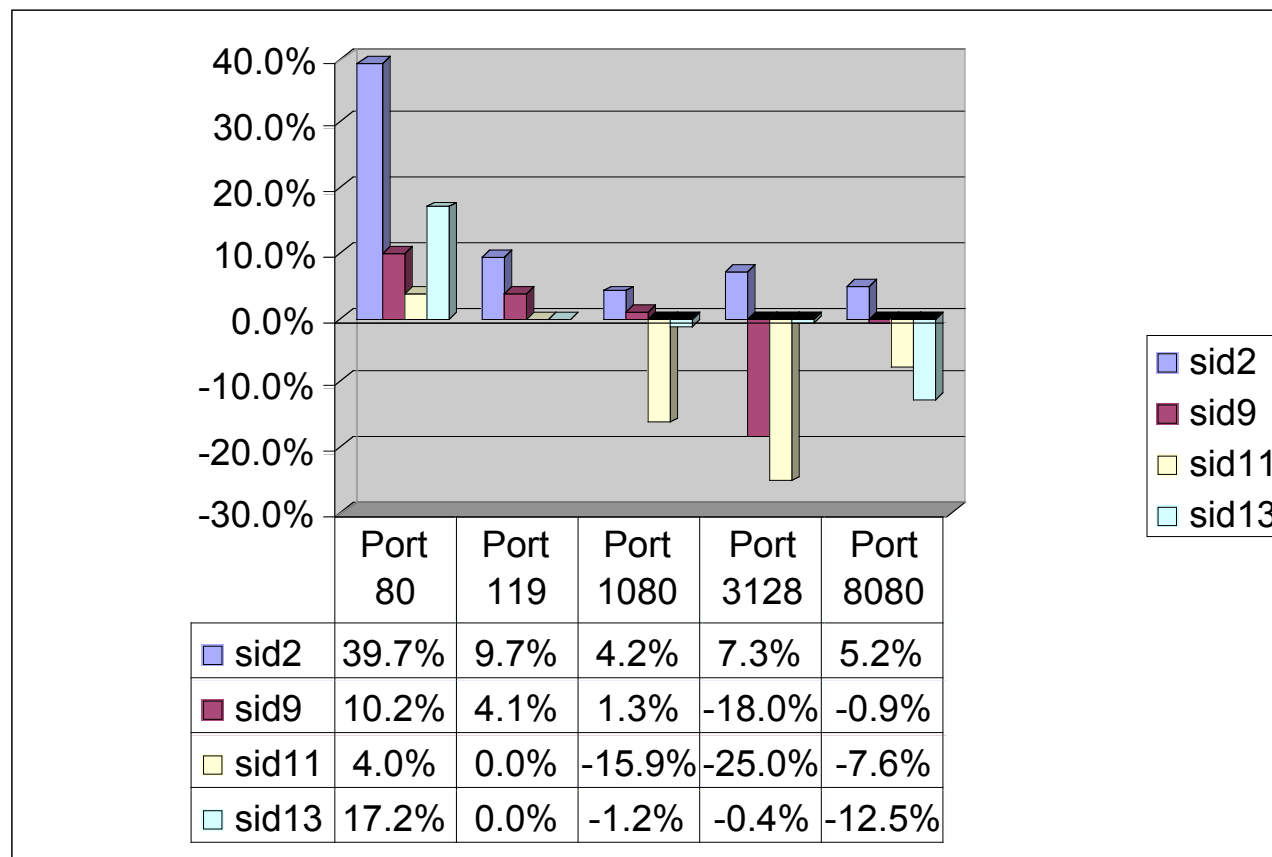
---

- No new attacks
  - Only matches known signatures
- Lack of context
  - Don't know what we are not seeing
- Non-standardized signature rule sets
  - No administrative control
- Missing Data
  - Uncertainty: Sensor failure vs. no intrusion attempts



## TCP Destination Port Changes

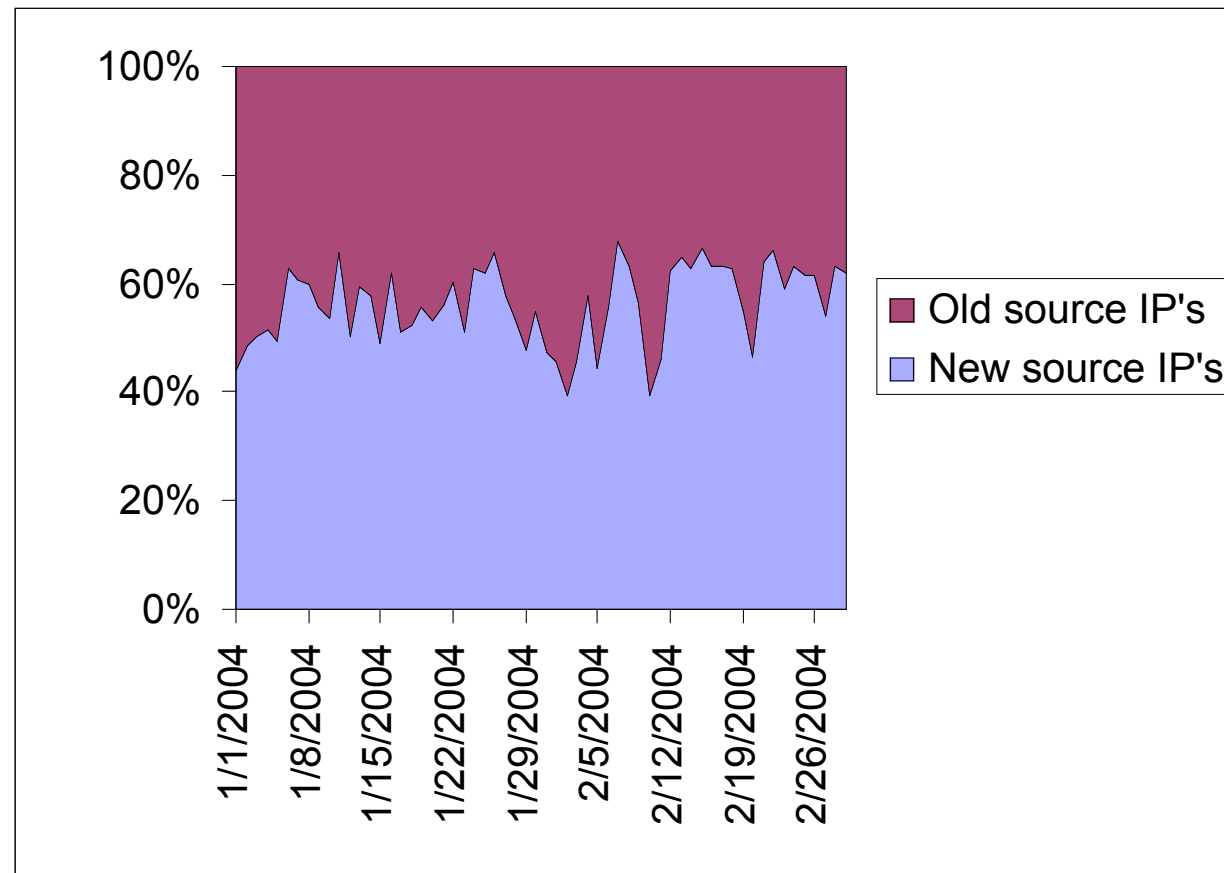
Comparison of port activity across organizations shows monthly trends.





## Share of New Source IP Addresses

Share of new daily source IP addresses stays fairly consistent.

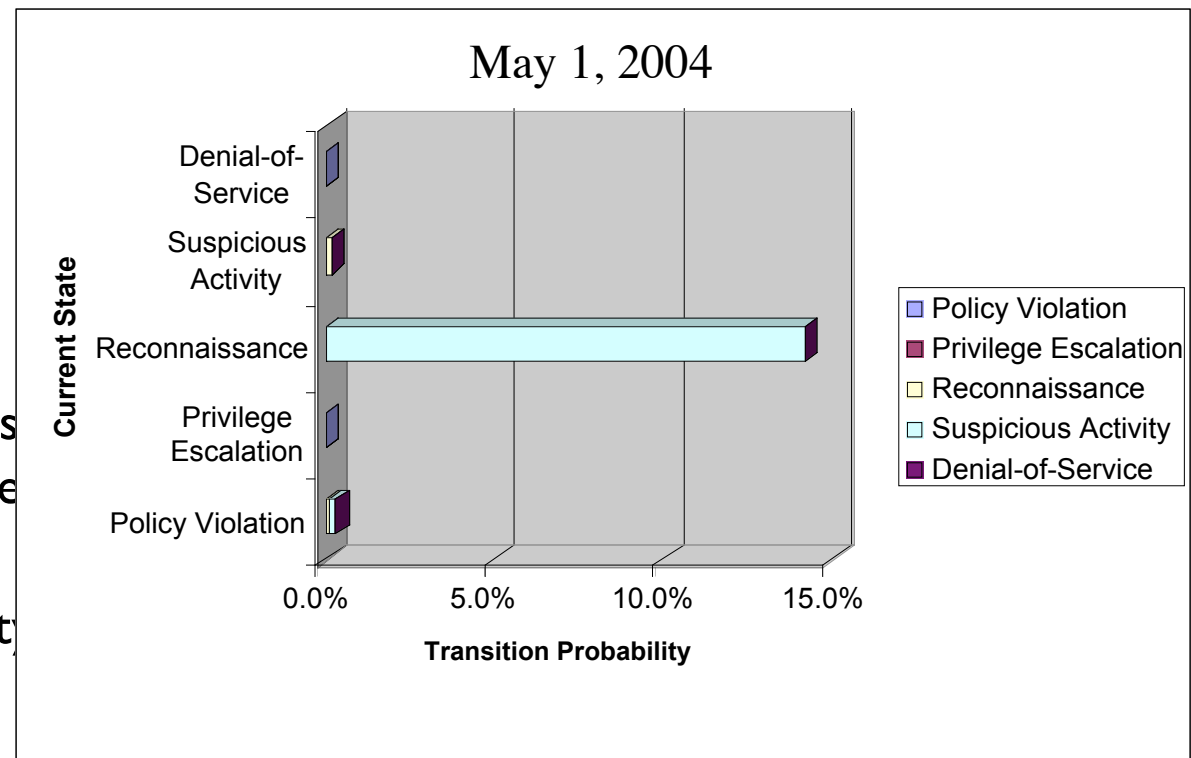




# Signature Class Transition

Transition probabilities highlight sequential patterns in data.

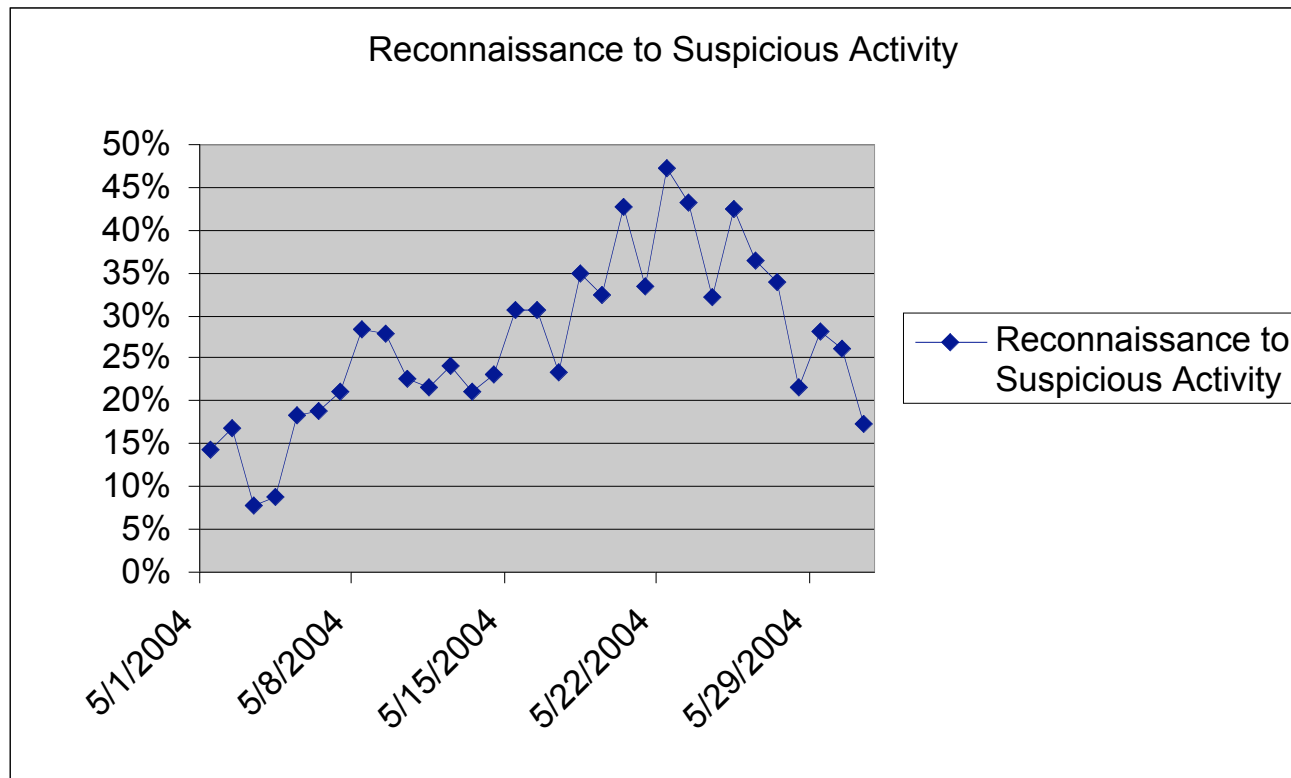
- **Current State**
  - Source IP records alert on Destination IP
- **Transition probability**
  - Percent chance for next class of alert recorded
- **Most source/dest combos involve only one signature class**
- **Small transition probability for**
  - Privilege Escalation





## Daily Transition Probabilities

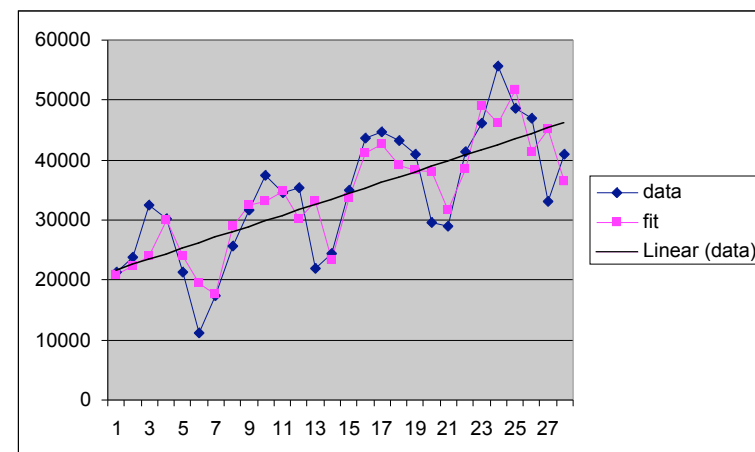
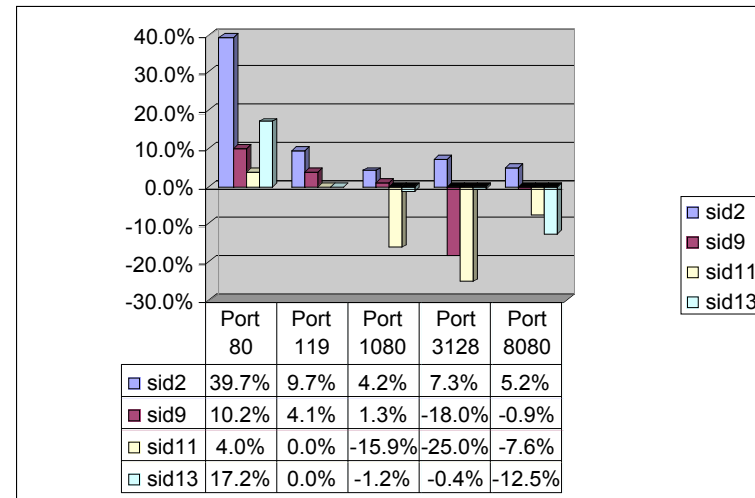
Transition probabilities can be monitored over time to identify consistent sequences.





# Trend Detection

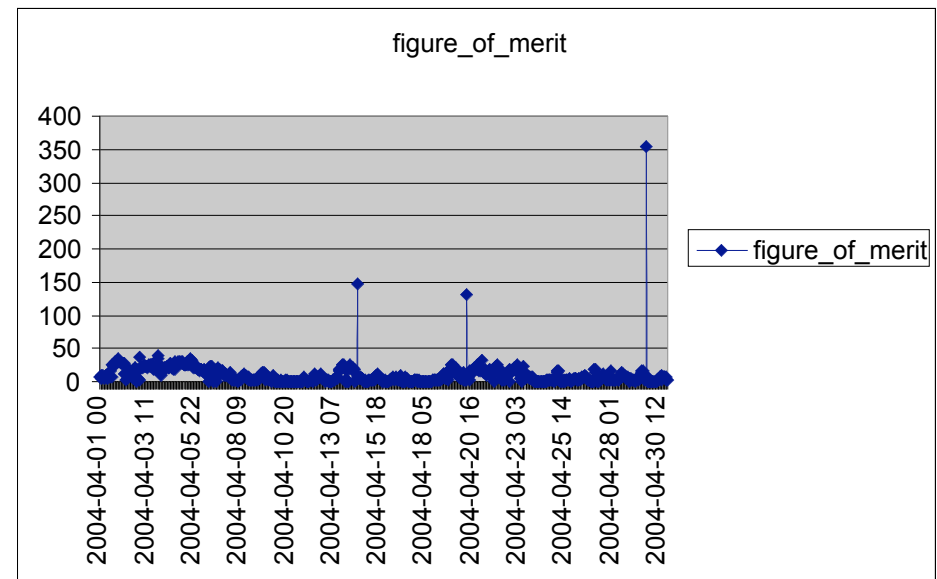
- Current month vs. previous month
  - Across organizations
  - % changes
- Time Series
  - Fit trend line
    - Arbitrary time period
  - Seasonal Components
  - Regression with ARMA errors





# Anomaly Detection

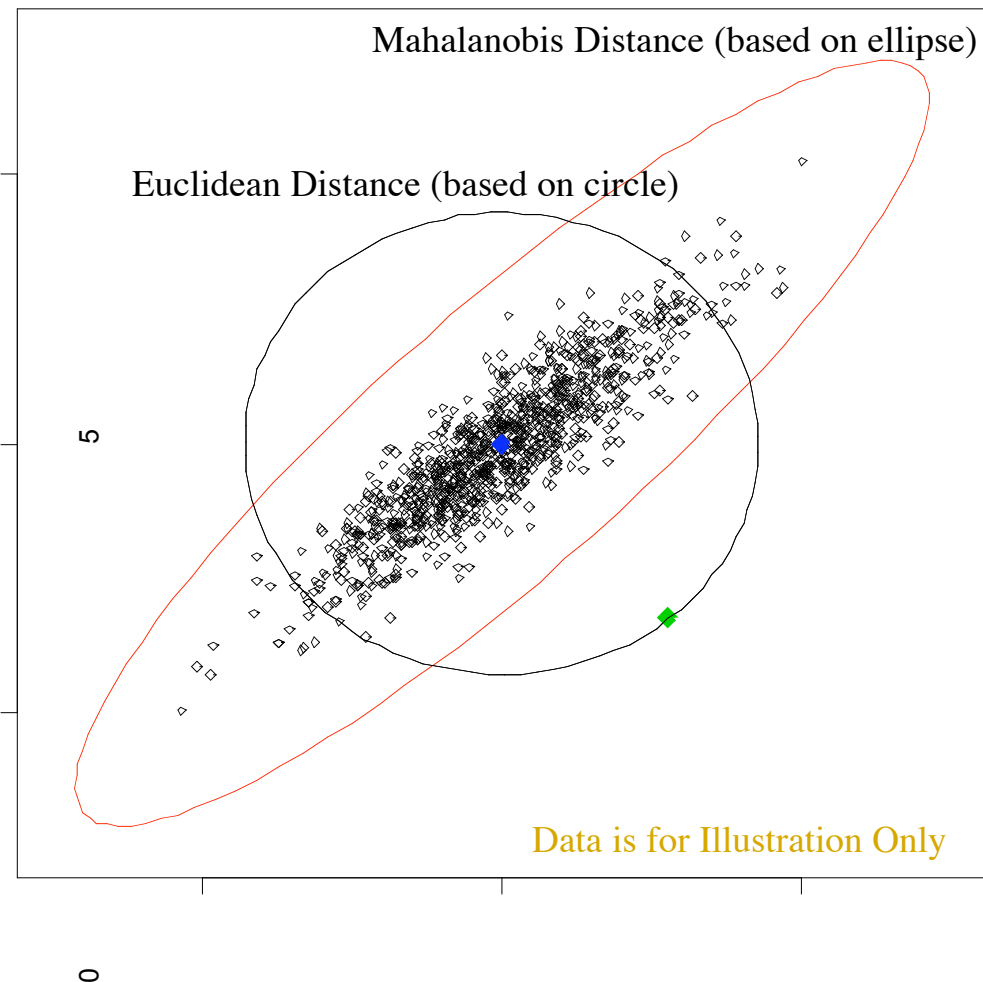
- Goal: Identify data points which deviate from overall pattern of data
- Our current implementation (Figure of Merit)
  - Evaluate hours
  - Record # alerts, # source IP addresses, # destination IP addresses, # signatures
- For each hour, we want measure of how deviant it was.





## Mahalanobis distance: 2D case

- Compute distance metric between each hour and the **average** hour
- When measuring **Euclidean** (**Mahalanobis**) Distance, all points along **circle** (**ellipse**) are same distance from the center
  - Points on larger circle/ellipse are greater distance from center
- Shape of the ellipse
  - Function of correlation between variables
- Generalizes to n dimensions (Ellipsoid)







## Analysis Roadmap

---

- Incorporate flow data
- Automating trend detection
  - Time series analysis
- Clustering
  - Group sources by similar activity patterns
    - Temporal correlation
    - Targeting similarities
    - Signature usage
  - Look for evidence of possible coordination

## Network Telescopes: The FloCon Files



 David Moore, Colleen Shannon  
{dmoore,cshannon}@caida.org  
[www.caida.org](http://www.caida.org) 

## Flocon Stream of Consciousness

- There are "reseachers" seriously interested in pieces of operational problems.
  - anomaly detection, early worm detection
  - flow aggregation, line-speed summarization
  - distributed data collection
  - modeling of "normal" traffic
- However, they can really use your help to understand the questions you currently ask and what you'd like to ask, but can't now.

 University California, San Diego – Department of Computer Science  
COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS 

## What is CAIDA?

- Cooperative Association for Internet Data Analysis
- Goals include measuring and understanding the global Internet.
- Develop measurement and analysis tools
- Collect and provide Internet data: topology, header traces, bandwidth testlab, network security, DNS
- Visualization of the network

 University California, San Diego – Department of Computer Science  
COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS 

## Current Project Areas

- Routing topology and behavior
- Passive monitoring and workload characterization
- Internet Measurement Data Catalog
- Bandwidth estimation
- Flow collection and efficient aggregation
- Security: DoS and Internet worms, syslog/SSH
- DNS performance and anomalies
- Visualization
- P2P traffic detection and modelling

 University California, San Diego – Department of Computer Science  
COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS 

## Tools

- CoralReef, NeTraMet, cflowd – packet, flows
- Walrus & Otter, libsea, PlotPaths - visualization
- NetGeo – IP to geography (mostly defunct)
- Skitter – large scale traceroute
- Graph::Chart.pm, GeoPlot.pm – plotting
- ASFinder.pm – IP to prefix/AS from routing table
- Beluga, GTrace – user-level traceroute viz
- dnstat, dnstop – passive DNS analysis
- DBHost, OWL – historical network meta-data (whois, DNS)
- Collaborations:
  - RRDTool, AutoFocus, PathRate/PathLoad



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## What is a "Network Telescope"?

- A way of seeing remote security events, without being there.
- Can see:
  - victims of certain kinds of denial-of-service attacks
  - hosts infected by random-spread worms
  - port and host scanning
  - misconfiguration



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Network Telescope

- Chunk of (globally) routed IP address space
- Little or no legitimate traffic (or easily filtered)
  - might be "holes" in a real production network
- Unexpected traffic arriving at the network telescope can imply remote network/security events
- Generally good for seeing explosions, not small events
- Depends on statistics/randomness working



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Amount of Telescope Data

- Currently collecting 30G/day of compressed data, and this is not including NetBios.
- Some "real-time" web reporting.
- Keep packet headers for a couple days, more summarized data longer, everything automatically rolled off to tape archive system.



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Flat File Compression

- Heard bzip2, gzip.
- We really like **lzop** for many things. It's close to gzip -1 size, but: faster, block-based, block checksums, ...
- Both lzop, gzip -1:
  - Allows packet capture to disk at higher data-rates.
  - Allows faster wall-clock analysis on datasets.
- bzip always slow: compressing and decompressing.



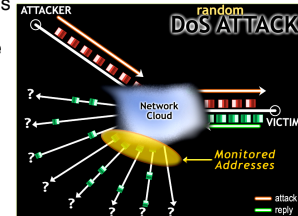
University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Network Telescope: Denial-of-Service Attacks

- Attacker floods the victim with requests using random spoofed source IP addresses
- Victim believes requests are legitimate and responds to each spoofed address
- With a /8 ("class A"), one can observe 1/256<sup>th</sup> of all victim responses to spoofed addresses



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Assumptions and Biases

- *Address uniformity*
  - Ingress filtering, reflectors, etc. cause us to **underestimate** number of attacks
  - Can bias rate estimation (can we test uniformity?)
- *Reliable delivery*
  - Packet losses, server overload & rate limiting cause us to **underestimate** attack rates/durations
- *Backscatter hypothesis*
  - Can be biased by purposeful unsolicited packets
    - Port scanning (minor factor at worst in practice)
  - Can we verify backscatter at multiple sites?



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Backscatter Hypothesis Busted?

- Not all TCP RST packets are DoS backscatter.
- Have seen a distributed scan using TCP RST packets spread over more than a month
  - "random" /25s (128 victim IPs) at a time, from a ~100 hosts, looking for a couple specific ports. TTL is not low. Seen at more sites than our /8.
- What were they trying to find? Current best guess, looking for differential ICMP error responses.

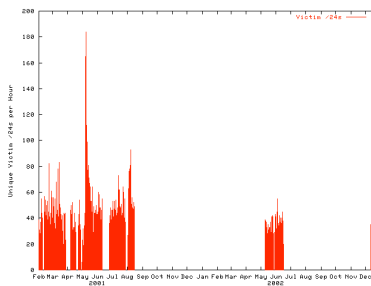


University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## DoS Attacks over time



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Our Telescope Data Analysis

- "Flow" based
  - Packets collected where possible, but most initial analysis is done with tools which work on flow-like aggregates.
- Eg, for backscatter
  - look at "outdegree" of victim IPs to telescope addresses



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## E.G. backscatter

- "Keys":
  - victimIP, protocols
- "Counters":
  - #pkts
  - #telescope IPs (also some distribution info)
  - #ports (also some distribution info) (for both src/dst)
  - are ports incrementing, decrementing (in little-endian byte order?)

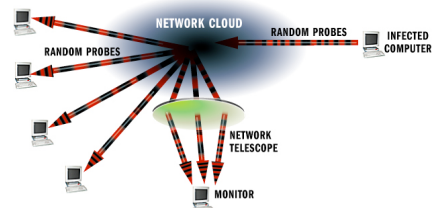


University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Network Telescope: Worm Attacks



- Infected host scans for other vulnerable hosts by randomly generating IP addresses
- We monitor 1/256<sup>th</sup> of all IPv4 addresses
- We see 1/256<sup>th</sup> of all worm traffic of worms (when no bias or bugs)



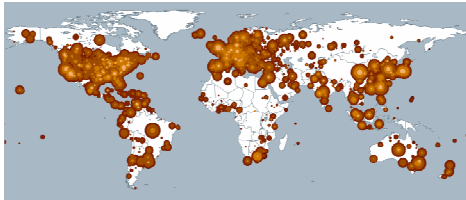
University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Internet Worm Attacks: Code-Red

(July 19, 2001)



- 360,000 hosts infected in *ten hours*, 2,000 new per minute at peak
- No effective patching response
- More than \$1.2 billion in economic damage in the first ten days
- Collateral damage: printers, routers, network traffic



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Response to August 1st CodeRed

- CodeRed was programmed to deactivate on July 20<sup>th</sup> and begin spreading again on August 1<sup>st</sup>
- By July 30<sup>th</sup> and 31<sup>st</sup>, more news coverage than you can shake a stick at:
  - FBI/NIPC press release
  - Local ABC, CBS, NBC, FOX, WB, UPN coverage in many areas
  - National coverage on ABC, CBS, NBC, CNN
  - Printed/online news had been covering it since the 19<sup>th</sup>
- “Everyone” knew it was coming back on the 1<sup>st</sup>
- Best case for human response: known exploit with a viable patch and a known start date



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Patching Survey

- How well did we respond to a best case scenario?
- Idea: randomly test subset of previously infected IP addresses to see if they have been patched or are still vulnerable
- 360,000 IP addresses in pool from initial July 19<sup>th</sup> infection
- 10,000 chosen randomly each day and surveyed between 9am and 5pm PDT

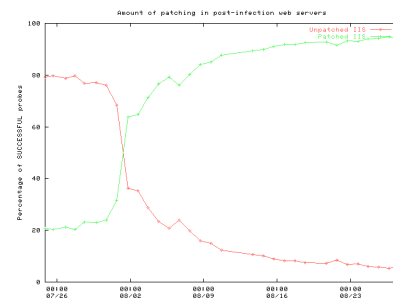


University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Patching Rate



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Dynamic IP Addresses

- How can we tell how when an IP address represents an infected **computer**?
- Resurgence of CodeRed on Aug 1st: Max of ~180,000 unique IPs seen in any 2 hour period, but more than 2 million across ~a week.
- This **DHCP effect** can produce skewed statistics for certain measures, especially over long time periods.
- Important to keep in mind if making big "bad lists".



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Dynamic IP Addresses

- For each /24, count:
  - total number of unique IP addresses seen ever
  - maximum number seen in 2 hour periods
- On plot:
  - x-axis is total number of unique addresses seen ever
  - y-axis is maximum number for a 2 hour period
  - the  $x = y$  (total = max) line shows /24s that had all their vulnerable hosts actively spreading in same 2 hour period, and those hosts didn't change IP addresses
  - the space far below and to the right of the  $x = y$  line (total  $\gg$  max) shows /24s that appear to have a lot of dynamic addresses
  - color of points represents density (3d histogram)

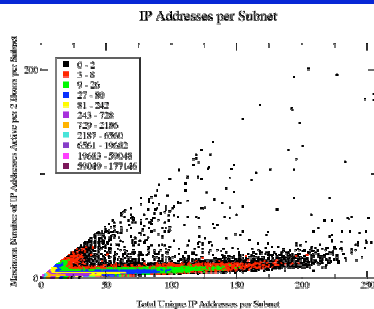


University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## DHCP Effect seen in /24s



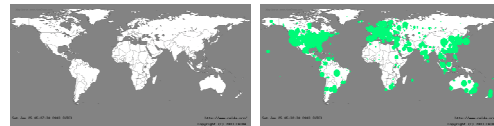
University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Internet Worm Attacks: Sapphire

(aka SQL Slammer) – Jan 24, 2003



- ~100,000 hosts infected in **ten minutes**
- Sent more than 55 million probes per second world wide
- Collateral damage: Bank of America ATMs, 911 disruptions, Continental Airlines cancelled flights
- Unstoppable; relatively benign to hosts



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Spread of the Witty Worm

March 19, 2004

- First wide-spread Internet worm with destructive payload  
writes 64k blocks to disk at random location, repeatedly
- Launched from a large set of ground-zero hosts  
>100 hosts
- Shortest interval from vulnerability disclosure to worm release  
1 day
- Witty infected firewall/security software  
i.e. proactive user base
- Spread quickly even with a small population  
~12,000 total hosts, 45 minutes to peak of infection

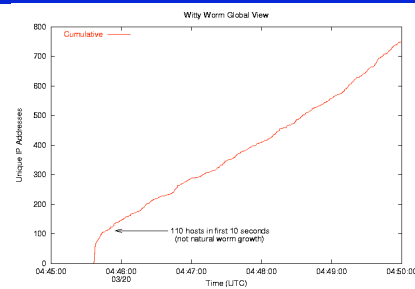


University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Early Growth of Witty

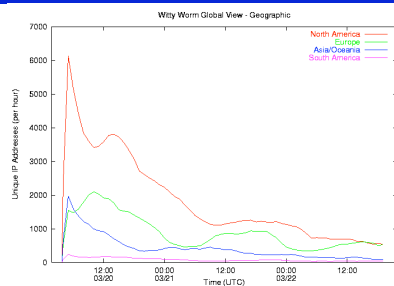


University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Geographic Spread of Witty



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Passive Vulnerability Fingerprinting

- Really good idea, helps finds new devices/services on the network.
- Minor (?) downside: miss new services running which aren't actually used.
- Recent major downside: Miss many *passive* devices in the network.
  - transparent caches, proxies, BlackIce Defender..., AMP boxes



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS





## Conclusions

- Don't really have conclusions, but it seems like there are some good community opportunities out there. I don't see any transit ISP security folks here, some of them are currently using netflow (or passive devices).
- Watch out for DHCP effect.
- Watch out for passive devices in network.
- Academics generally don't understand operational needs, make lists.



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Related CAIDA/UCSD Papers

- Inferring Internet Denial-of-Service Activity [MSV01]
  - David Moore, Stefan Savage, Geoff Voelker
  - <http://www.caida.org/outreach/papers/2001/BackScatter/>
- Code-Red: A Case Study on the spread and victims of an Internet Worm [MSB02]
  - David Moore, Colleen Shannon, Jeffrey Brown
  - <http://www.caida.org/outreach/papers/2002/codered/>
- Internet Quarantine: Requirements for Containing Self-Propagating Code [MSVS03]
  - David Moore, Colleen Shannon, Geoff Voelker, Stefan Savage
  - <http://www.caida.org/outreach/papers/2003/quarantine/>
- The Spread of the Sapphire/Slammer Worm [MPS03]
  - David Moore, Vern Paxson, Stefan Savage, Colleen Shannon, Stuart Staniford, Nicholas Weaver
  - <http://www.caida.org/outreach/papers/2003/sapphire/>



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Additional CAIDA/UCSD Information

- Code-Red v1, Code-Red v2, CodeRedII, Nimda
  - <http://www.caida.org/analysis/security/code-red/>
- Code-Red v2 In-depth analysis
  - [http://www.caida.org/analysis/security/code-red/coderedv2\\_analysis.xml](http://www.caida.org/analysis/security/code-red/coderedv2_analysis.xml)
- Spread of the Sapphire/SQL Slammer Worm
  - <http://www.caida.org/analysis/security/sapphire/>
- Network telescopes
  - <http://www.caida.org/analysis/security/telescope/>



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Using your own telescope: Effects of Size

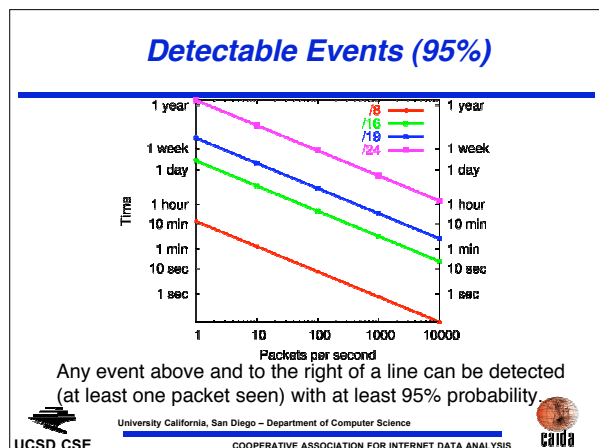
- Larger telescopes are able to detect events that generate fewer packets, either because of short duration or low sending rate.
- Larger telescopes have better accuracy at determining the start and end times of an event.
- Using CIDR / notation on next few slides:
  - /8 = old class-A size, 16 million IP addresses
  - /16 = old class-B size, 65536 IP addresses



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



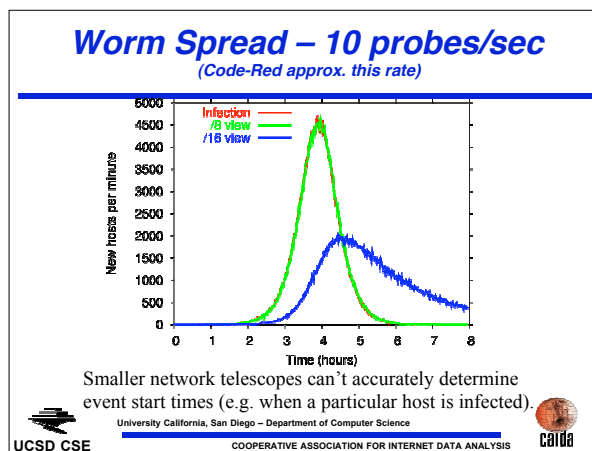
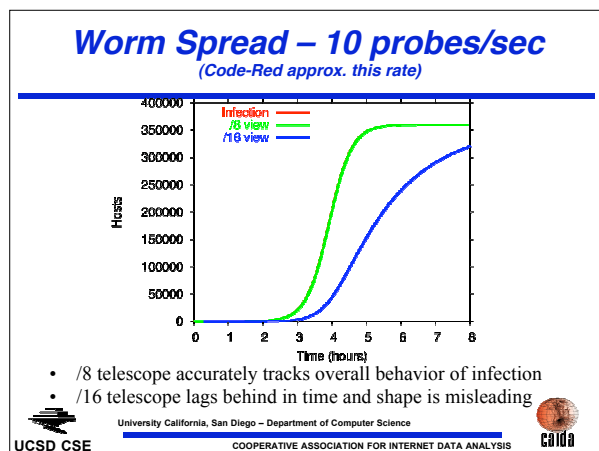


### Detection Times - 10 pps events

(Code-Red approx. this rate)

Detection probability:	5%	50%	95%
/8	1.3 sec	18 sec	1.3 min
/14	1.4 min	19 min	1.4 hour
/15	3 min	38 min	2.7 hour
/16	6 min	1.3 hour	5.5 hour
/19	45 min	10 hour	1.8 day
/24	24 hours	14 day	58 day

UCSD CSE University California, San Diego – Department of Computer Science COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS caida



## Organizational Telescopes

- Small telescopes may not be useful for observing external events
- However, setting up an internal facing telescope may help quickly identify internal problems
- With an internal facing telescope you can have /5 or better



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Why have an internal telescope?

- Quickly detect internal machines infected with worms, certain kinds of misconfigurations, and potentially hacked machines.
- Capture data for hosts connecting to unallocated IP address space by:
  - if you use BGP (default-free) to all providers, you can point a default route at a monitor box
  - enable flow collection on your edge routers
  - announce a couple unallocated networks, but be careful if they ever get allocated by IANA (least desirable)



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Extending it

- Combine a telescope watching traffic to unallocated IP addresses with monitoring all outbound traffic
  - you may notice anomalous behavior like a spam relay
  - verify that your firewall seems to be doing what you think
- Watch all *inbound* ICMP error messages, in particular HOST/NETWORK UNREACHABLE
  - evidence of scanning behavior
  - may show external connectivity & performance problems before users pick up the telephone



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Tools to use

- Flow data (Cisco NetFlow, Juniper cflow, others):
  - FlowScan: <http://net.doit.wisc.edu/~plonka/FlowScan>
- Packet data
  - CoralReef report generator: <http://www.caida.org/tools/>
- Either
  - AutoFocus: <http://jal.ucsd.edu/AutoFocus/>
- Not an exhaustive list ☺



University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## AutoFocus example

- Sapphire/SQL Slammer worm
- Find worm port & proto automatically

Source IP	Destination IP	Protocol	Source Port	Destination Port	bytes	Unspectacularity(%)
"	"	6	highports	highports	827M	77.7
"	"	17	highports	1434	10.5G	112.6
"	152.249.0.0/16	"	"	"	604M	100
138.0.0.0/9	"	"	highports	"	3.66G	99.4
138.0.0.0/10	"	"	highports	"	3.68G	99.9
138.54.3.58	"	17	3341	1434	2.14G	672.5
138.54.11.4	"	17	7062	1434	950M	1551.3
152.249.56.0/22	"	"	highports	highports	723M	103.4
152.249.191.120	"	17	1959	1434	1.78G	810.0
152.249.191.121	96.0.0.0/8	17	1531	1434	645M	39523.7
152.249.210.3	"	17	4315	1434	2.36G	609.5
152.249.254.152	"	17	3787	1434	1.53G	941.8

UCSD CSE

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



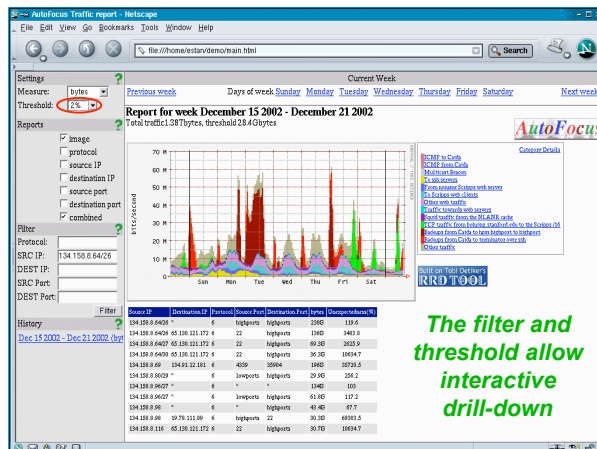
## AutoFocus example

- Sapphire/SQL Slammer worm
- Can identify infected hosts

Source IP	Destination IP	Protocol	Source Port	Destination Port	bytes	Unspectacularity(%)
"	"	6	highports	highports	827M	77.7
"	"	17	highports	1434	10.5G	112.6
"	152.249.0.0/16	"	"	"	604M	100
138.0.0.0/9	"	"	highports	"	3.66G	99.4
138.0.0.0/10	"	"	highports	"	3.68G	99.9
138.54.3.58	"	17	3341	1434	2.14G	672.5
138.54.11.4	"	17	7062	1434	950M	1551.3
152.249.56.0/22	"	"	highports	highports	723M	103.4
152.249.191.120	"	17	1959	1434	1.78G	810.0
152.249.191.121	96.0.0.0/8	17	1531	1434	645M	39523.7
152.249.210.3	"	17	4315	1434	2.36G	609.5
152.249.254.152	"	17	3787	1434	1.53G	941.8

UCSD CSE

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS



## Conclusions

- Network telescopes provide insight into non-local network events
- Larger telescopes better capture the behavior of events and can see smaller events
- Build your own internal telescope – it's fun AND easy.

UCSD CSE

University California, San Diego – Department of Computer Science

COOPERATIVE ASSOCIATION FOR INTERNET DATA ANALYSIS





# **Flow Data Analysis in SWITCH / ETH Zurich Project DDoSVax**

Arno Wagner

wagner@tik.ee.ethz.ch

Communication Systems Laboratory

Swiss Federal Institute of Technology Zurich (ETH Zurich)



# Talk Outline

---

- The Dataset
- Flow Data Usage by SWITCH
- Offline Analysis Examples
- Traffic Amount vs. Unique Addresses
- Analysis Tools
- Performance questions

# The DDoSVax Dataset



Project URL:

<http://www.tik.ee.ethz.ch/~ddosvax/>

- NetFlow v5 (converted from V7 by SWITCH)
- About 60.000.000 flows/hour
- Weekday: About 200k internal and 800k external IPs
- Unsampled
- Stored in full since March 2003



# Flow Data Usage by SWITCH



Independently done by SWITCH on NetFlow data

- Accounting and load monitoring (aggregated)
- SWITCH-CERT: Short-term forensics (reduced)
  - Single fast computer with hardware RAID-5
  - No compression
  - Sorted into minute (?) intervals
  - Fast search with regular expressions
  - Several weeks online
  - No (?) long term storage

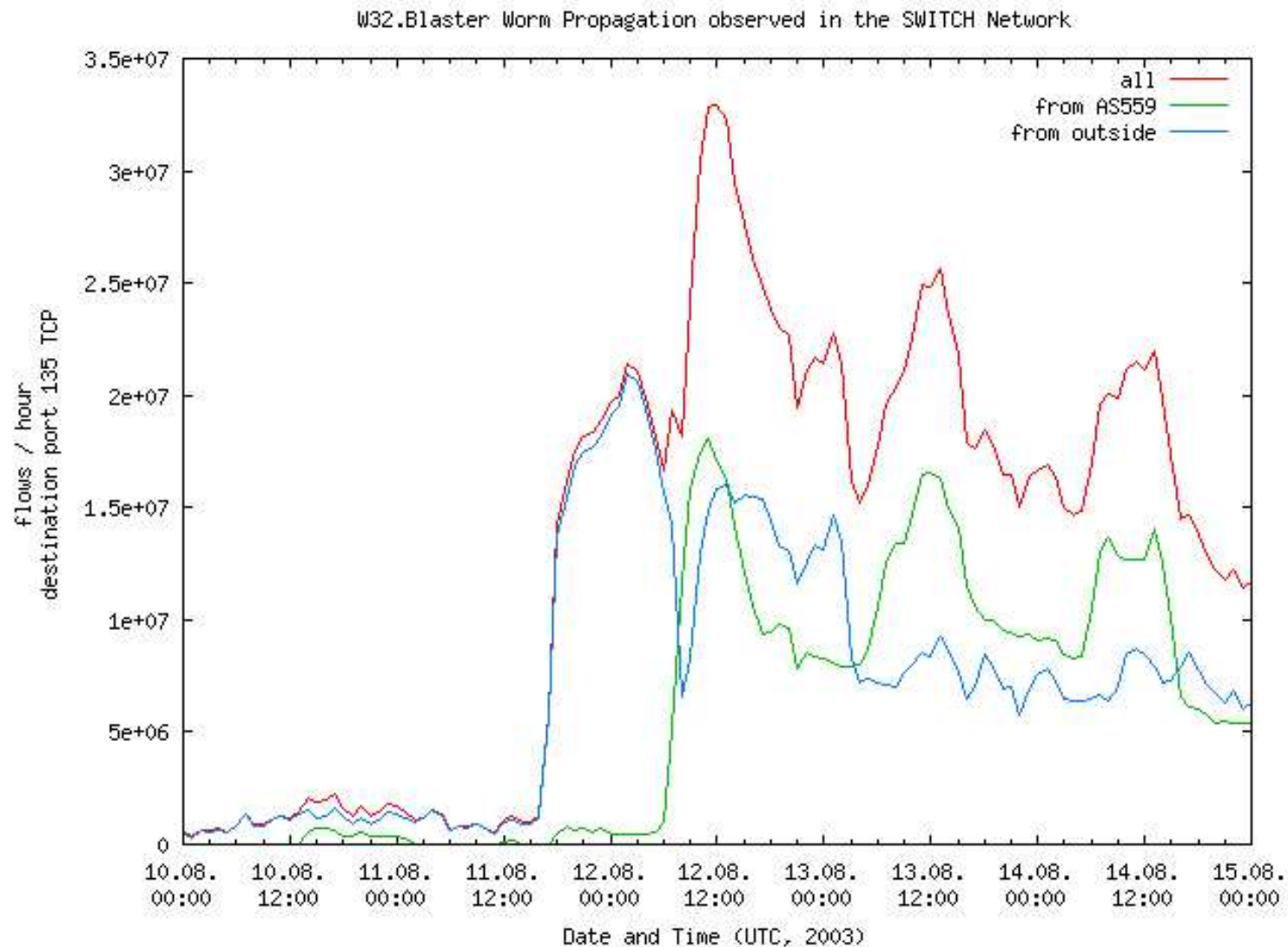




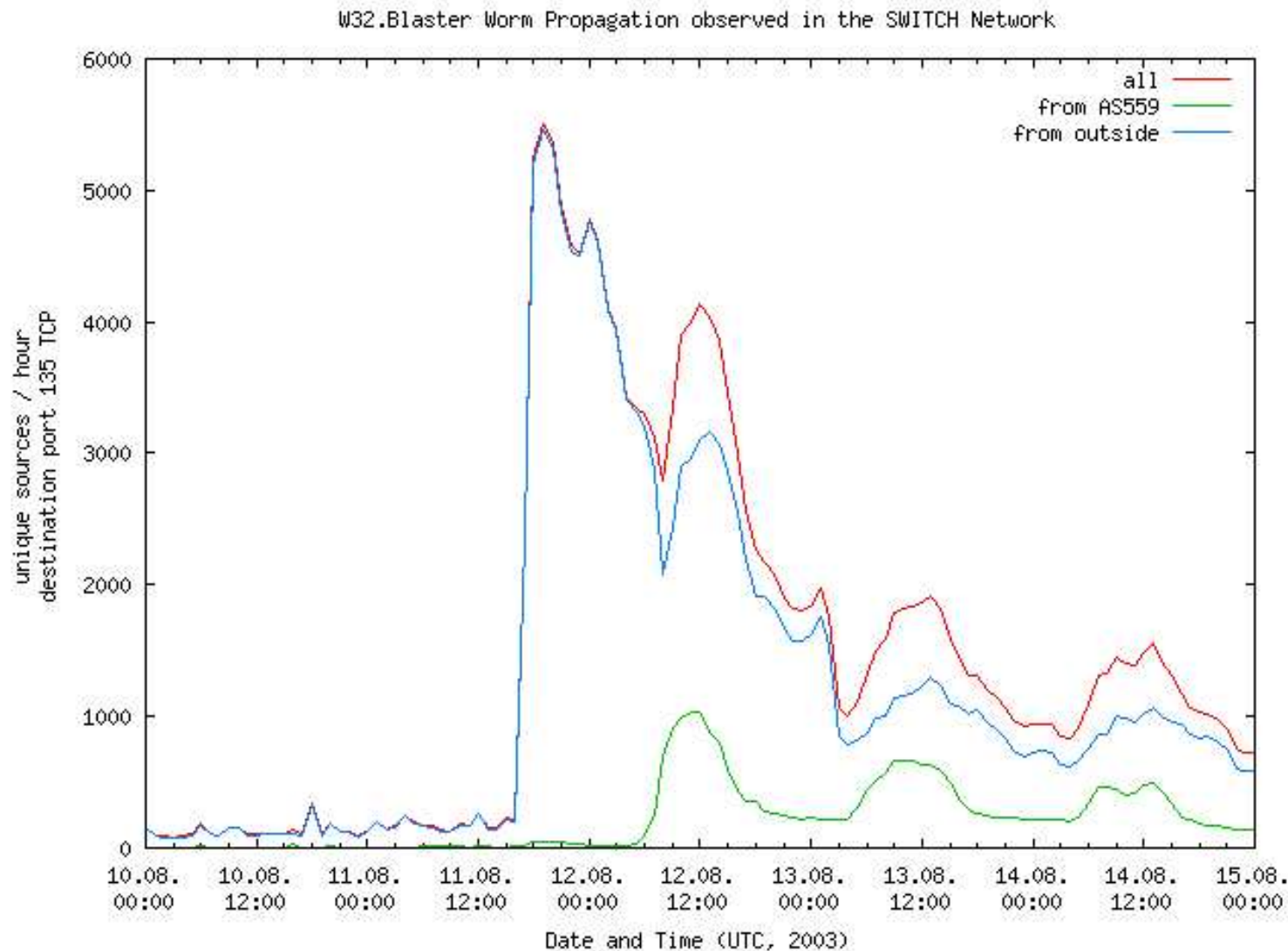
# Offline Analysis

- E.g. for network/email worms
- Customised tools for some analyses
  - Single hour / prototyping: netflow\_to\_text and Perl
  - Days...weeks: From C-template
- Also other things: P2P, IRC, ...

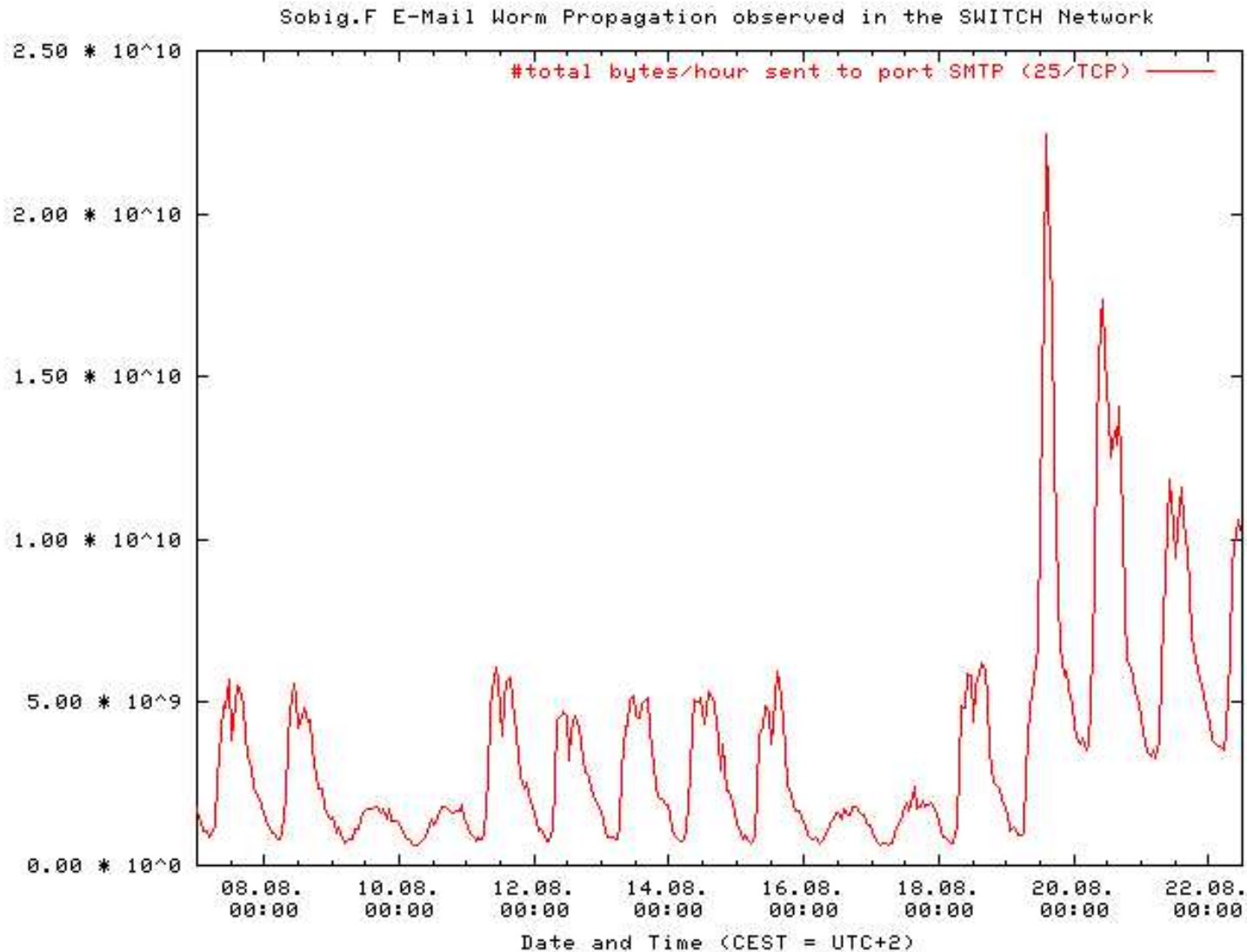
# Example: Blaster - Flows



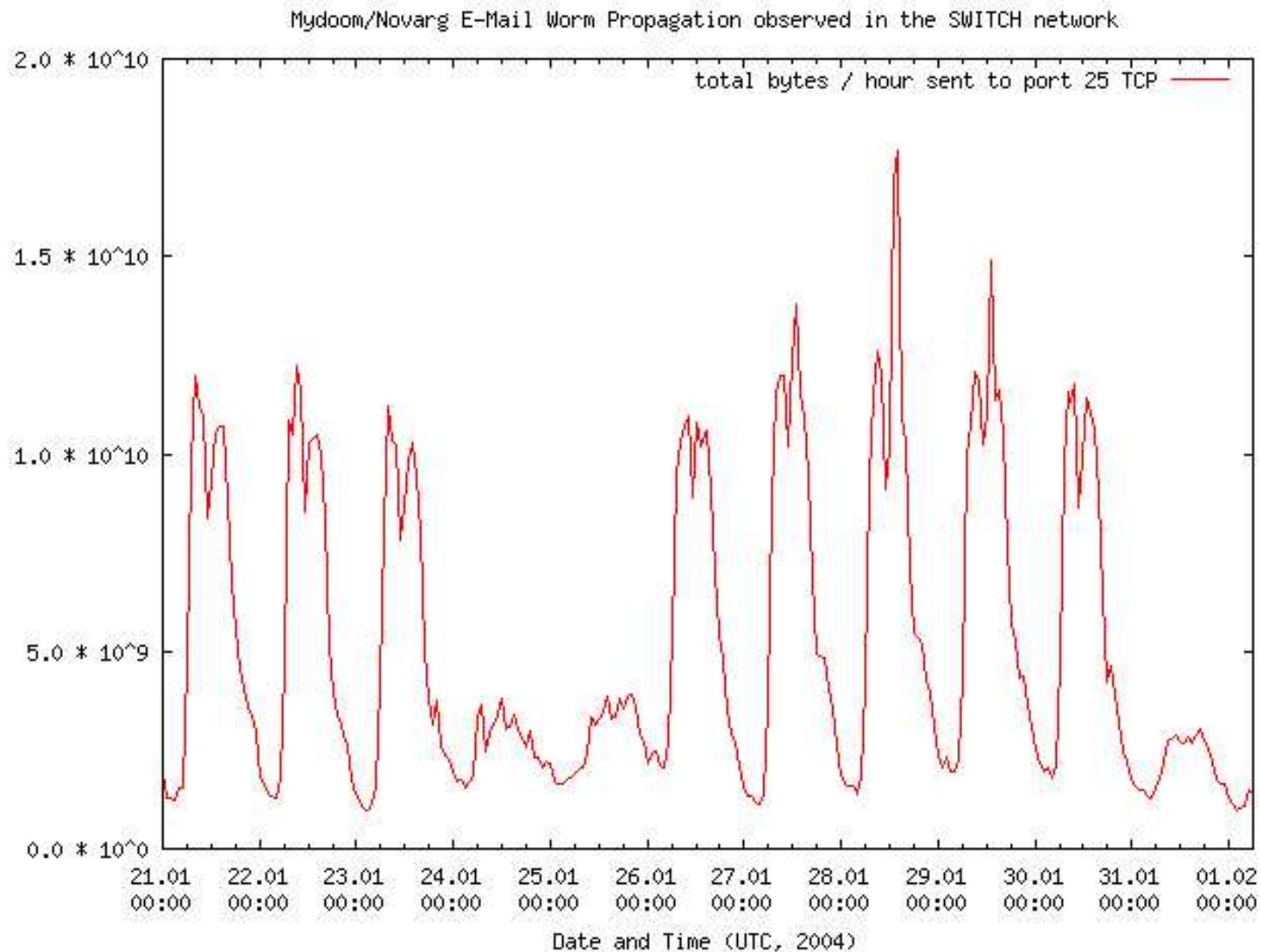
# Example: Blaster - Unique Sources



# Example: Sobig



# Example: MyDoom



# Traffic vs. Unique Sources



## Traffic:

- Easy to do
- Works reasonably well
- Sensitive to data generation problems
- Sensitive to observed network

## Unique Sources:

- More complicated, more robust
- Weakly dependent on observed network
- Allows to get global picture



# Analysis-tools: Scripting



"netflow\_to\_text"

- Takes one data file, outputs one line
- Well suited as "grep"/Perl input

Example:

```
TCP pr 111.131.210.8 si 1111.136.200.121  
di 1264 sp 135 dp 48 le 1 pk  
12:59:51.965 st 12:59:51.965 en 0.000 du
```



# Analysis-tools: C



"Iterator template"

- Iterates over all records in a set of files
- Preprocesses timestamps, etc.
- Reading of input files encapsulated





# Performance Issues

- 5-10 minutes / hour of data bunzip2
- I/O limit at 10 cluster nodes reading from one NFS partition
- Memory limitations



**Carnegie Mellon  
Software Engineering Institute**

**CERT**  
Situational  
Awareness

# AirCERT: Building a Framework for Cross-Administrative Domain Data Sharing

Roman Danyliw <rd@cert.org>

FloCon 2004: Complementary Architecture Panel

CERT® Network Situational Awareness Group  
Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

*The CERT Network Situational Awareness Group is part of the Software Engineering Institute. The Software Engineering Institute is sponsored by the U.S. Department of Defense.*





## Background

---

- Form situational awareness for the SEI, its sponsors, and the Internet community
  - Big picture view of threats
- Constraints
  - Situational awareness can only be formed with data from many organizations – all data is governed by the constraints of its owners
  - Must provide a reasonable value-proposition for data sharing
  - Strict hierarchies in data sharing will not scale
  - Solutions must be built with open and transparent architectures



# Analytical Concerns

---

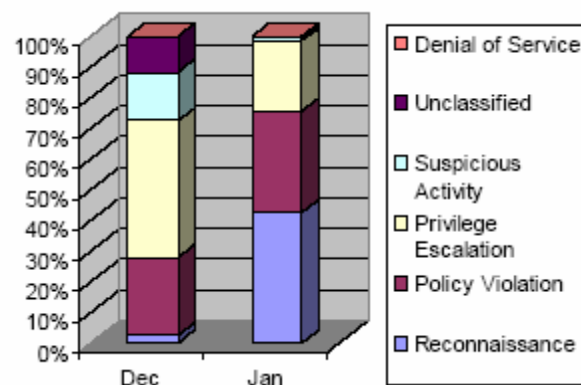
*Focus on merging and analyzing data from multiple view points*

- Distinguish between targeted, localized, and Internet-wide activity
  - Widely targeted services
  - Clusters of attacks
    - Passive detection of new tools
  - Attack techniques *de-jour*
  - Attack sources
- Historical trending
  - Enable forward estimation of expected intruder activity of a site

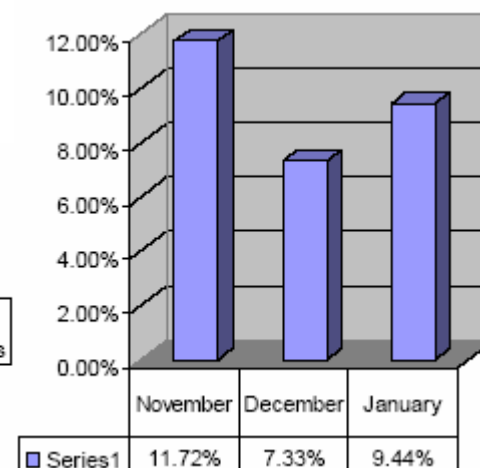
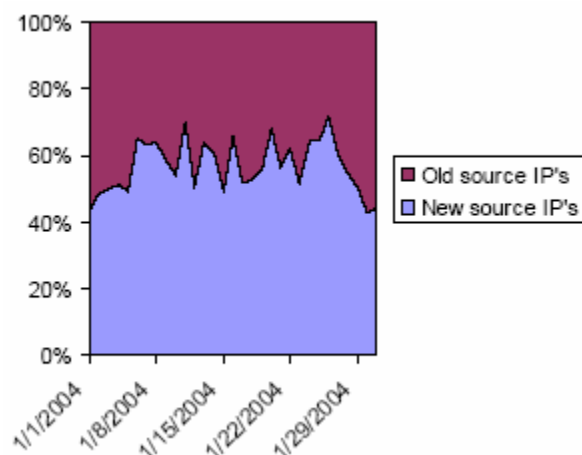


## Current Results

- Generating “Top 10” lists and volumetric measures based on
  - *Packet/Flow features*: IP addresses, ports, protocols, signature, etc.
  - *Context*: timing, vulnerability, country, net-blocks, etc.



Share Source IP addresses  
Targeting multiple  
organization





## Implementation

---

- <http://aircert.sourceforge.net>
- Gather data from existing security solutions already deployed
  - Partner with security operations in the federal civilian community and in academia
- Write “glue” to integrate, convert, analyze, and share the data across organizations
- Provide analytical results back to participants and sponsors



## Synthesized Data

---

- Categorization
  - SIM/SEMs (e.g., ArcSight)
  - IDS (e.g., Snort)
- Discovery
  - Flow data (e.g., argus)
- Refinement
  - Network topology information
  - IT/data data sharing policies
- Context
  - Vulnerability (e.g., CERT/CC KB)
  - Artifacts (e.g., CERT/CC AC)



## Collection Infrastructure

---

- Provides infrastructure to *automatically* extract relevant information from existing instrumentation
  - If human intervention is required, sharing is too expensive
- Wrote “normalizers” to handle the reformatting and semantic transformation of the data
  - Too many vendor to write one-off tools for each
  - Write transformation engine that understands the underlying data-store: text files, RDBMS, etc.





## Sharing Infrastructure: Collection

---

- The key to facilitating data sharing across organizations is
  - Making it seamless – no human interaction
  - Ensuring policy compliance
- All “normalizers”, “publishers”, and the underlying storage architecture have a notion that all data has an owner
  - Dissemination respects site’s local policy
  - Sanitization of sensitive data
  - Tagging of all data with a source identifier



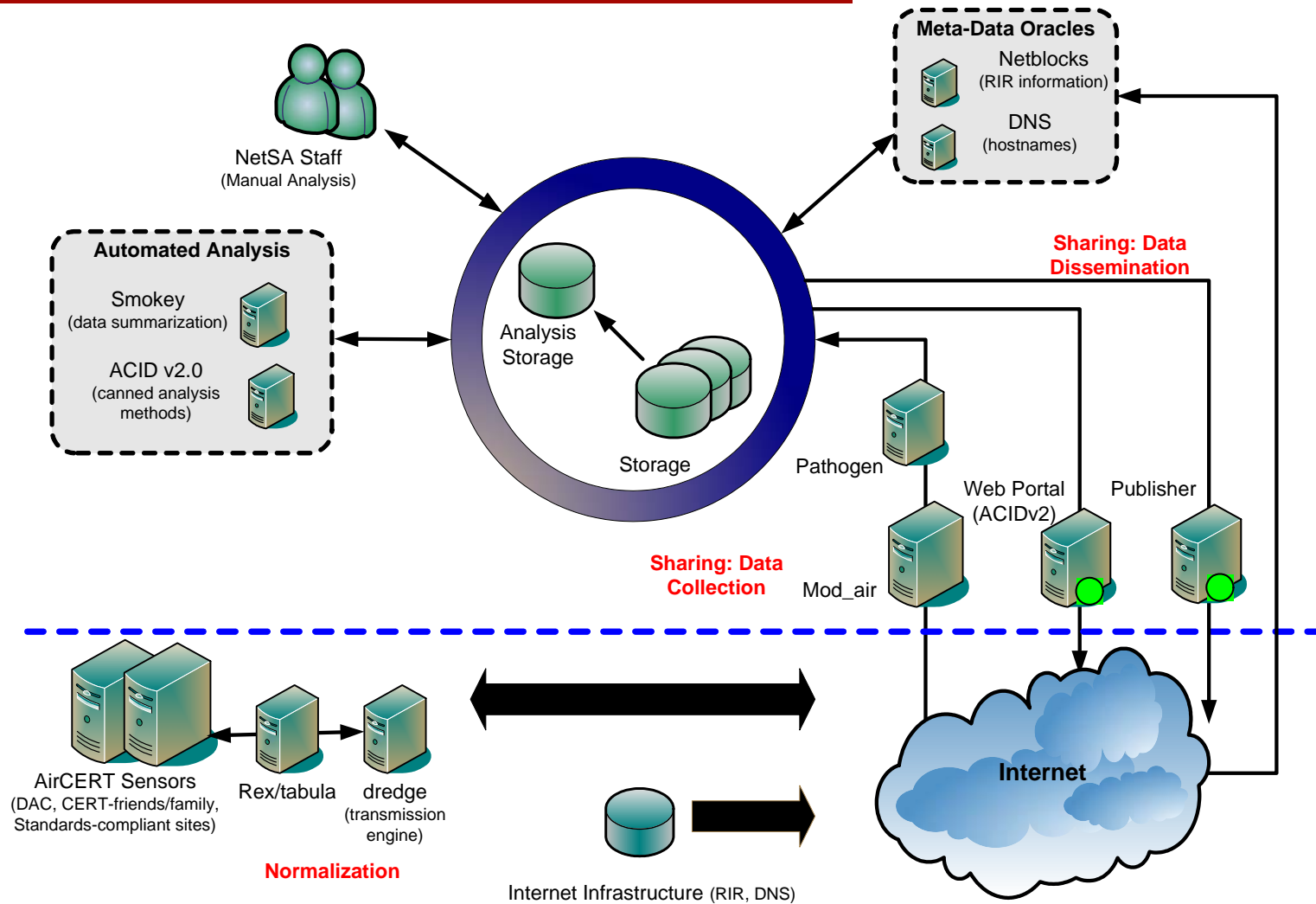
## Sharing Infrastructure: Dissemination

---

- Sharing data with us, is no different than data with others
- Tailor channel for the audience
  - Web-portal for pre-digested snapshot
  - Export bulk-data in a machine-readable format (e.g., XML, RSS)

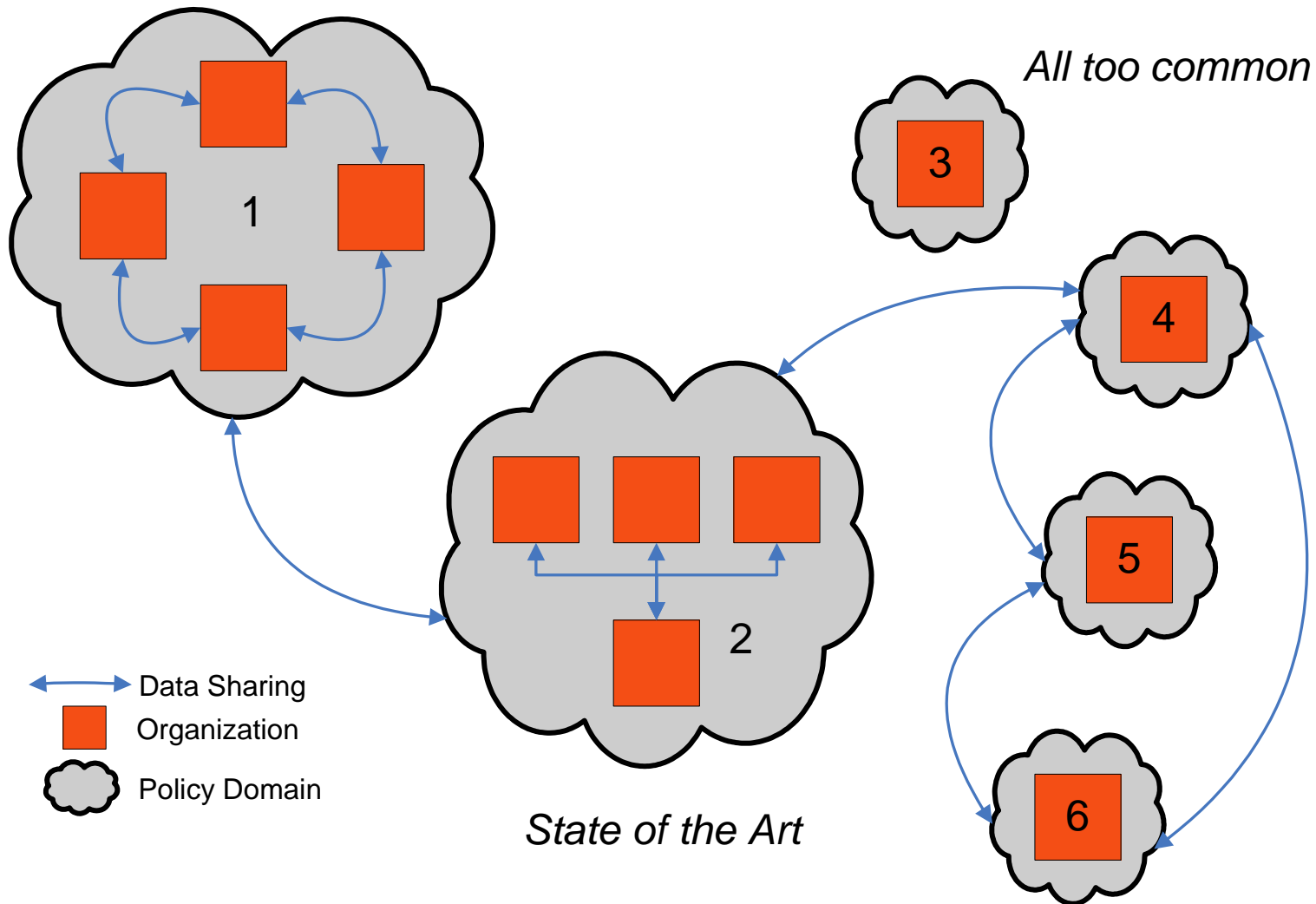


# Architecture





# Big Picture Architecture





## Challenges and Solutions

---

- Many different formats used by the SEM and IDS products
  - Support standards efforts: IDMEF, IODEF, IPFIX, PSAMP
  - Storage-specific normalization tools
- Normalizing signatures across IDS products
  - Using CVE and custom classification taxonomies
- Analyzing the correct signature set
  - Use only explicit malicious activity
  - Filtering out policy violations and poorly written signatures
  - Use the correct tool for the task
    - Deploy non-IDS sensors next to the IDS
- Data loops
  - “Checksums” in the meta-data of the data stream



## Challenges and Solutions

---

(2)

- Need both push and pull model, while supporting varied levels of automation
  - Unified presentation engine (ACIDv2)
  - Publisher for bulk-data transfer



## Ongoing Work

---

- Intelligent end-points that summarize instead of sending all data
- Automated ways to overlay the context provided by vulnerability and artifact information
- Continued support for standards work
- Improved attention focusing techniques for flow data-to-IDS and vice versa
- Improved approaches for integrating the analytical products into operations



# **NetFlow Data Capturing and Processing at SWITCH and ETH Zurich**

Arno Wagner

wagner@tik.ee.ethz.ch

Communication Systems Laboratory

Swiss Federal Institute of Technology Zurich (ETH Zurich)





# Talk Outline

---



- The DDoSVax Project
- The SWITCH Network
- NetFlow Data Capturing Infrastructure
- Long-Term Storage
- Computing infrastructure
- Infrastructure Cost
- Remarks and Lessons Learned
- Online Processing Framework: UPFrame
- Conclusion



# The DDoSVax Project



`http://www.tik.ee.ethz.ch/~ddosvax/`

- Collaboration between SWITCH ([www.switch.ch](http://www.switch.ch)) and ETH Zurich ([www.ethz.ch](http://www.ethz.ch))
- Aim (long-term): Analysis and countermeasures for DDoS-Attacks and Internet Worms
- Start: Begin of 2003
- Funded by SWITCH and the Swiss National Science Foundation



# SWITCH

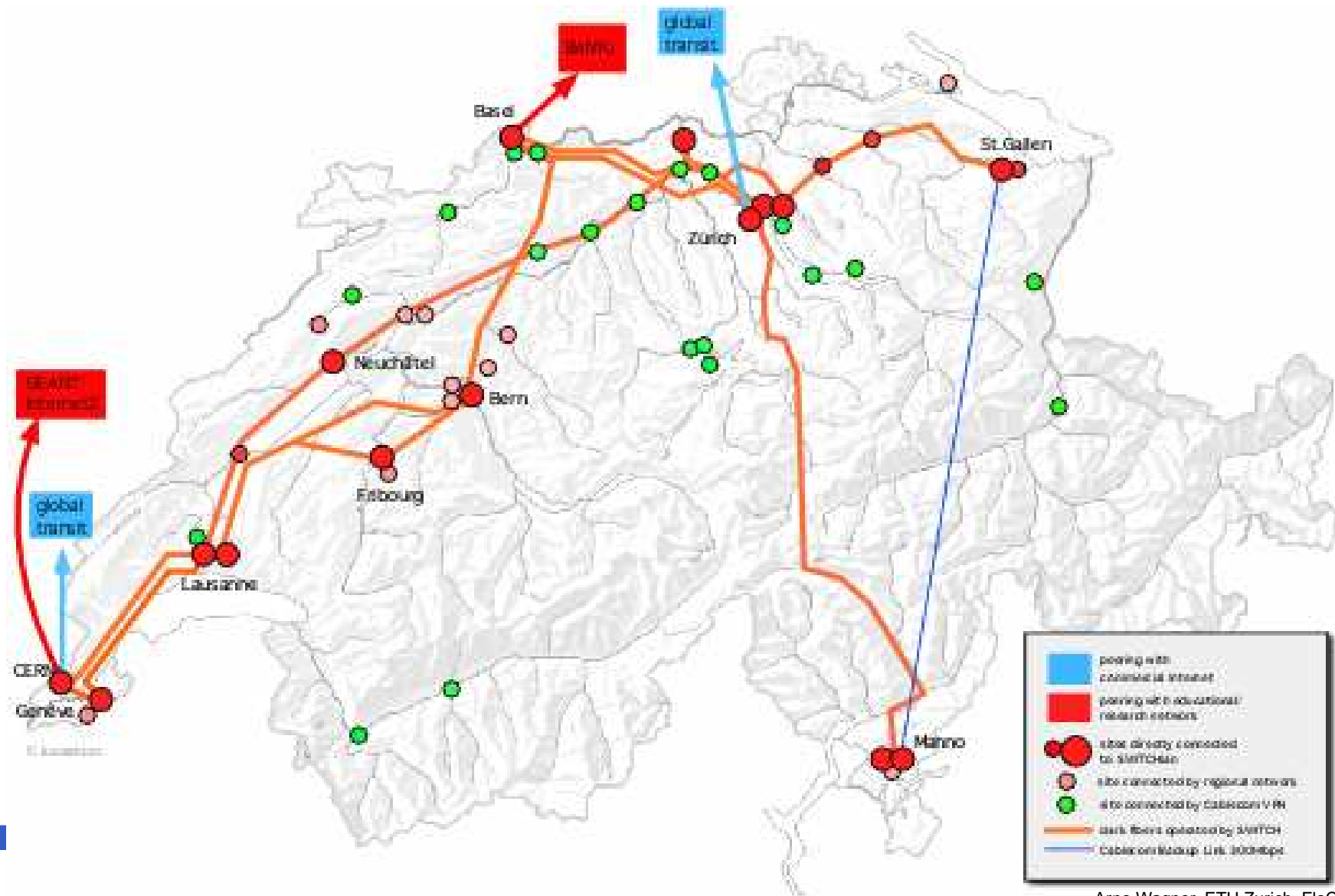


## The Swiss Academic And Research Network

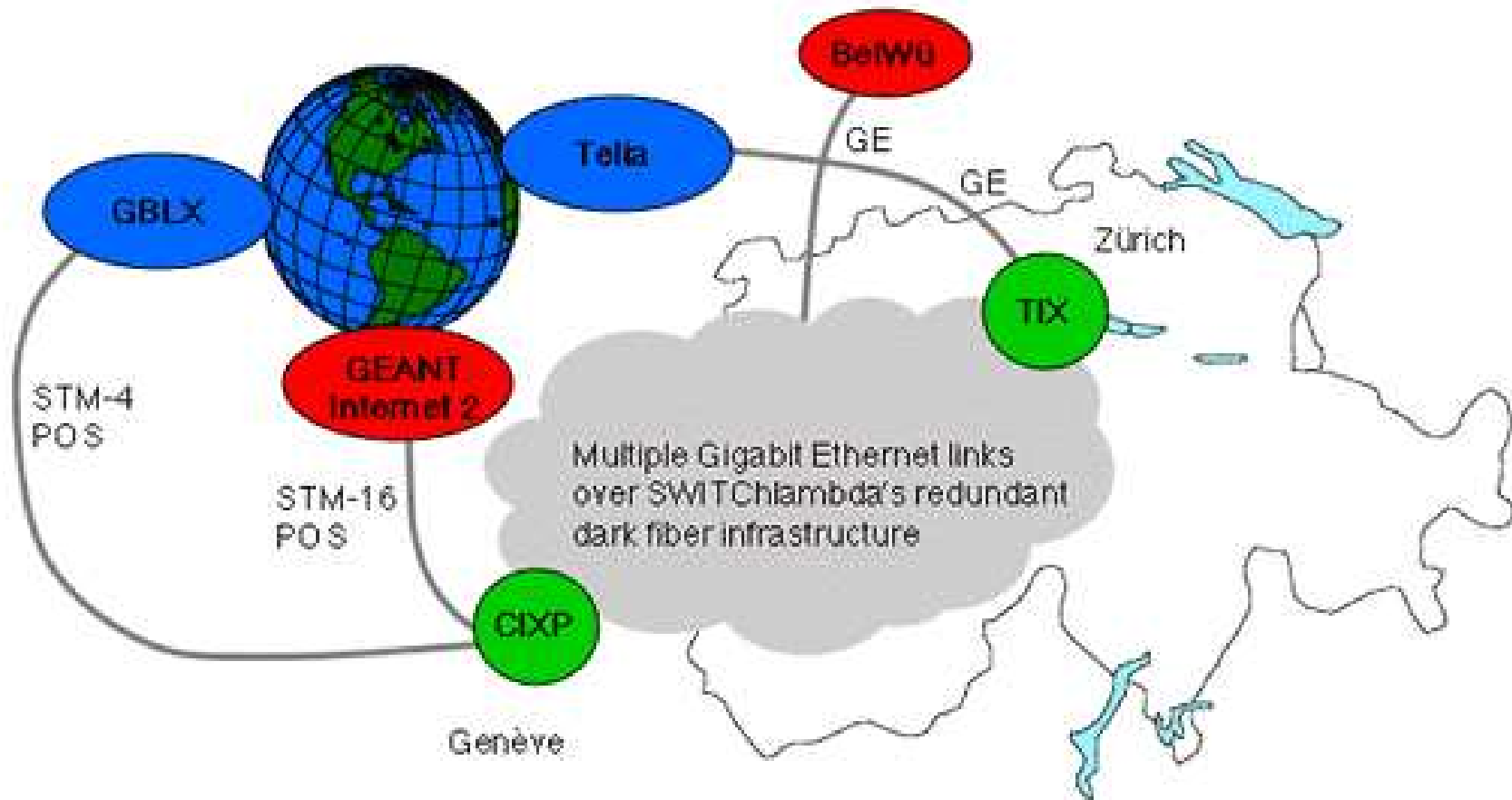
- .ch Registrar
- Links most (all?) Swiss Universities
- Connected to CERN
- Carried around 5% of all Swiss Internet traffic in 2003
- Around 60.000.000 flows/hour
- Around 300GB traffic/hour






# The SWITCH Network



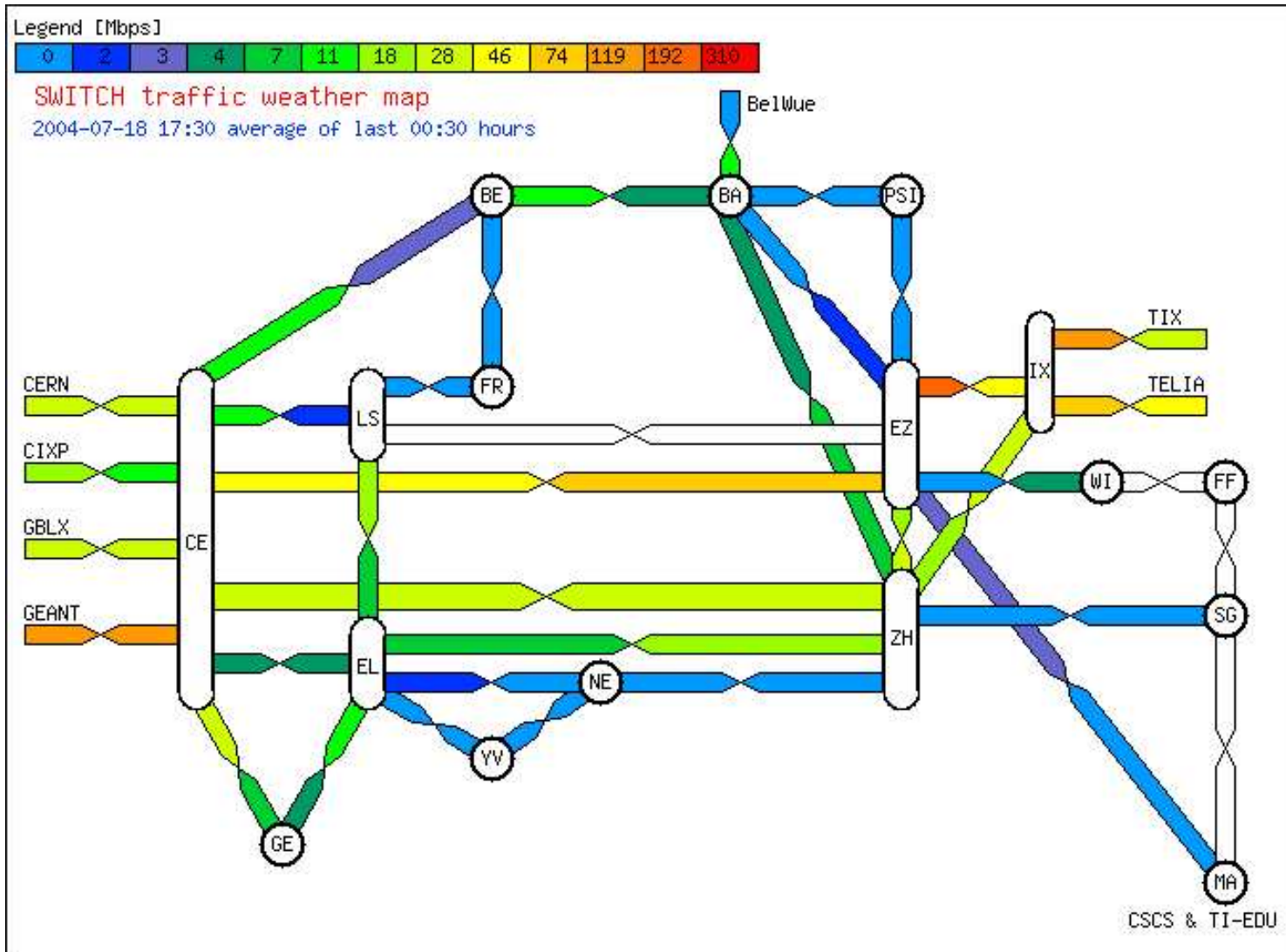
# SWITCH Peerings



-  Global transit by international carriers
-  Private peering with international research networks

-  Public Internet eXchange with bilateral peerings

# SWITCH Traffic Map



# SWITCH Routers



(Don't ask me for specifics...)

- swiCE2, swiCE3, swiX1: Cisco 7600 OSR with Supervisor 720
- swiBA2: Cisco 7600 OSR with Supervisor 2
- Cards: 8/16/48 GbE, 10GbE
- OSM POS OC-48c
- OSM POS 2\*OC-12c
- OSM 4\*Gigabit Ethernet



# NetFlow Data Usage at SWITCH



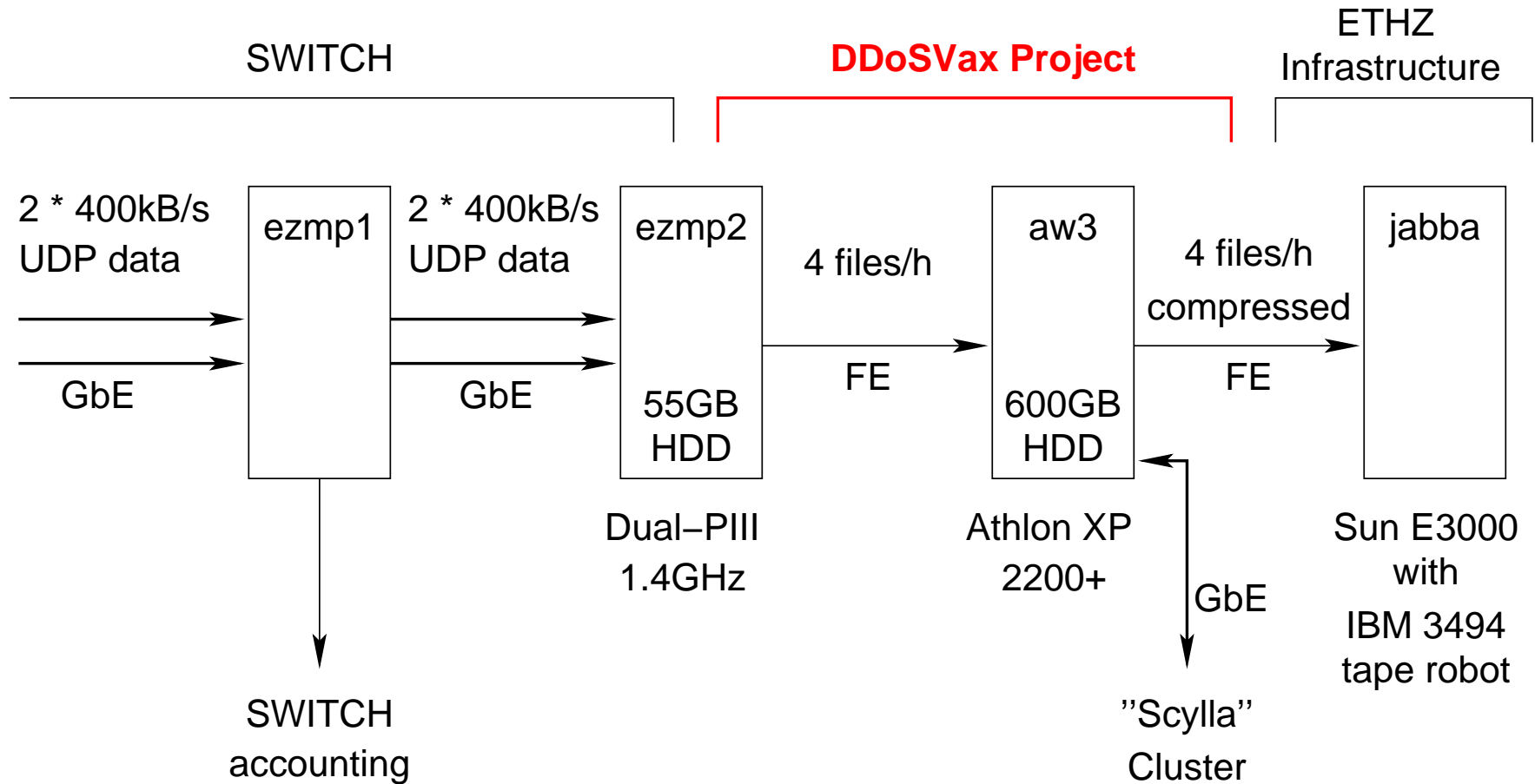
- Accounting
- Network load monitoring
- SWITCH-CERT, forensics
- DDoSVax (with ETH Zurich)

Transport: Over the normal network





# NetFlow Data Flow



# NetFlow Capturing

- One Perl-script per stream
- Data in one hour files
- Timestamps and src-IP in "stat" file

Critical: Linux socket buffers:

- Default: 64kB/128kB max.
- Maximal possible: 16MB
- We use 2MB (app-configured)
- 32 bit Linux: May scale up to 5MB/s per stream

# Capturing Redundancy

- Worker / Supervisor (both demons)
- Super-Supervisor (cron job)  
For restart on reboot or supervisor crash
- Space for 10-15 hours of data

No hardware redundancy

# Data Transfer to ETHZ



- Cron job, every 2 hours
- Single Perl script
- Transfer: scp (no compression, RC4)
- Remote deletion: ssh

No compression on ezmp2. (Some other Software running there)

Bzip2 compression on ezmp2 would be possible!



# Long-Term Storage Format



Full data since March 2003

Bzip2 compressed raw NetFlow V5 in one-hour files

- We need most data and precise timestamps
- We don't know what to throw away
- We have the space
- Preprocessing for specific work still possible

Latency: 5-10 minutes / hour of data



# Computing Infrastructure



## The "Scylla" Cluster Servers:

- aw3: Athlon XP 2200+, 600GB RAID5, GbE
- aw4: Dual Athlon MP 2800+, 800GB RAID5, GbE
- aw5: Athlon XP 2800+, 800GB RAID5, GbE

## Nodes:

- 22 \* Athlon XP 2800+, 120GB, GbE



# Infrastructure Cost Today

Speaker: 1 MYr = 175.000 CHF = 142.000 USD

⇒ 1MM = 12.000 USD, 1MD = 640 USD

Hardware and full installation:

- aw3 (capturing): 1600 USD + 2 MD
- aw4 (dual CPU server): 2500 USD + 3 MD
- Cluster: 24.000 USD + 1MM
- Maintenance: 1-2 MD/month

Hidden cost: Computer room, network infrastructure, software development

Scalability: Add 2\*200GB HDD to each node

⇒ 8TB additional at 6000 USD

# Lessons learned



Most important: KISS!

- Use scripting wherever possible
- Worker and Supervisor pairs are simpler  
⇒ "crash" as error recovery model
- Cron as basic reliable execution service
- Email for notification: Do rate-limiting
- File-copy: Interlock and age check
- ssh, scp password-less (user key)
- Nothing needs to run as "root"!





# Remarks on Software

- Linux is stable enough
- Linux is fast enough
- Linux Software RAID1/5 works well
- XFS has issues with Software RAID
- Perl is suitable for demons
- Python is suitable for demons

# Remarks on Hardware



PC hardware works well, but:

- Get good quality components (PSUs!)
- Get good cooling (HDDs/CPU's)
- Do SMART monitoring
- Do regular complete surface scans
- Have cold spares handy
- ...



# Remarks on Linux Clusters

- Rackmount vs. "normal"
- Cooling / Power needs planning
- Gigabit Ethernet "star" topology is nice
- KVM not for all nodes needed
- FAI (Fully Automatic Installation) for installation
- Local Debian mirror
  - ⇒ 10 Min for complete reinstallation
- No global connectivity for the nodes
- Private addresses for the nodes

# UPFrame

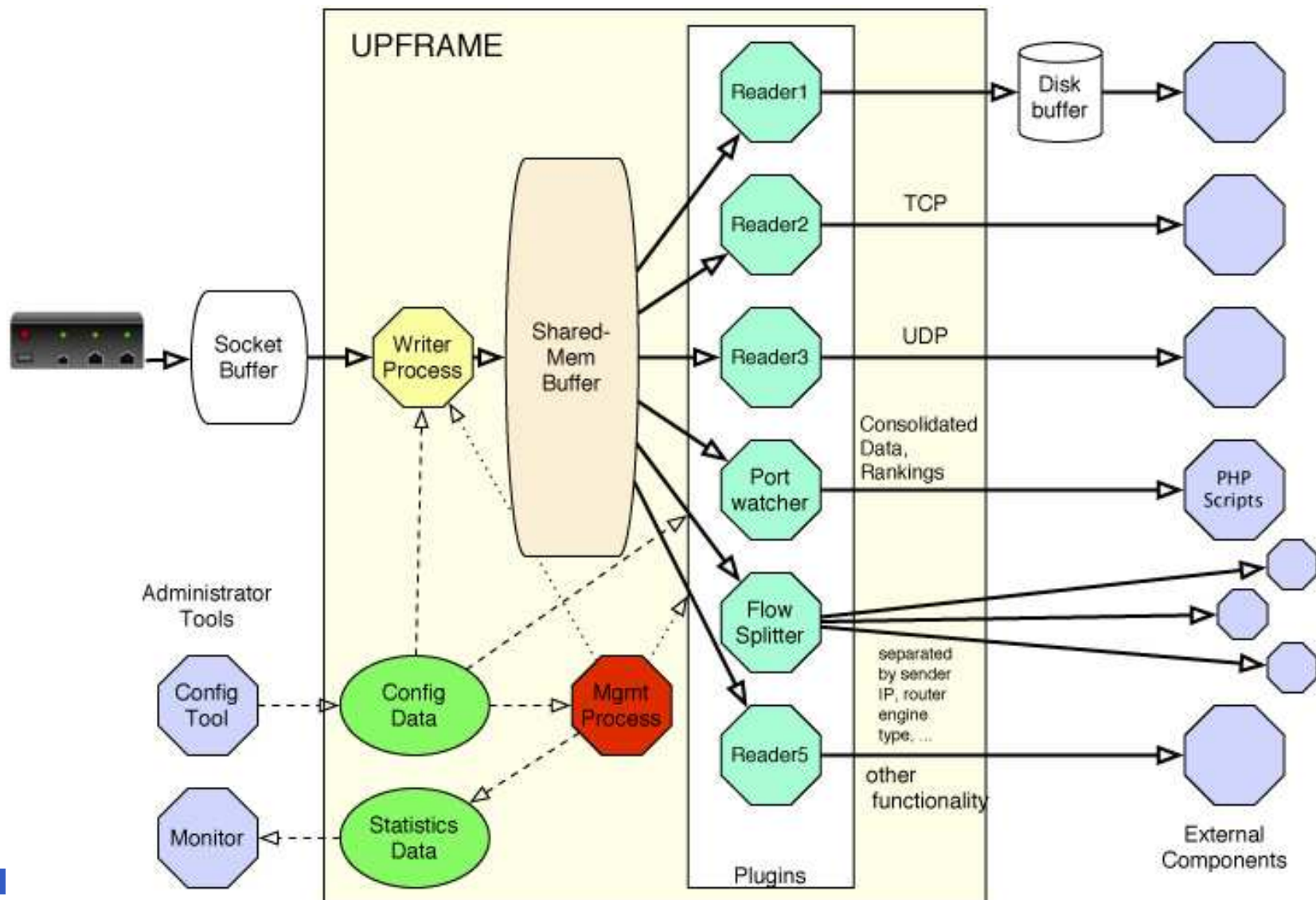


<http://www.tik.ee.ethz.ch/~ddosvax/upframe/>

- UDP plugin framework
- E.g. for online analysis of NetFlow data
- Can be used as traffic-shaper
- Robust: For experimental plugins



# UPFrame Structure



# Conclusion



- SWITCH is large enough and small enough
- No special hardware / software needed for capturing
- Long-term storage is unproblematic
- Linux can be used in the whole infrastructure
- Online processing is more difficult
- Simplicity and Reliability are the main issues
- ...





# Security at Line Speed with NetFlows

**William Yurcik**

*NCSA Security Research*

*National Center for Supercomputing Applications (NCSA)*

*University of Illinois at Urbana-Champaign*

*FloCon 2004*

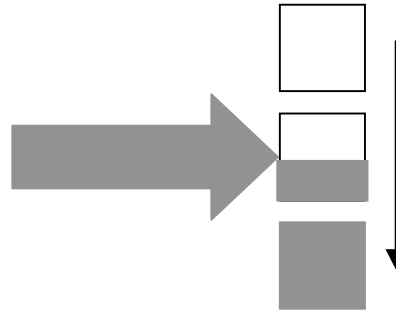
# Level of Observation

Data Source	Description	Advantage	Disadvantage
Packet	lowest level of granularity; all raw packets with all fields intact	most detailed data and statistics especially protocols; easiest to obtain	unscalable; protocol signaling needs to be decoded
NetFlows	IPs/ports/protocols/ Timestamps/data?	scalable for catching all traffic; multiple sources, uniform field formats	maybe no data field; context must be inferred
IDS	alerts of different formats	scalable; tunable	resource-intensive; misses; FPs
Load Levels	aggregate utilization levels that can be broken down to IP, protocol, port	high volume attacks (DOS, traffic); capacity planning; availability from routers & sniffers	details about SD pairs; no direction; low volume events obscured

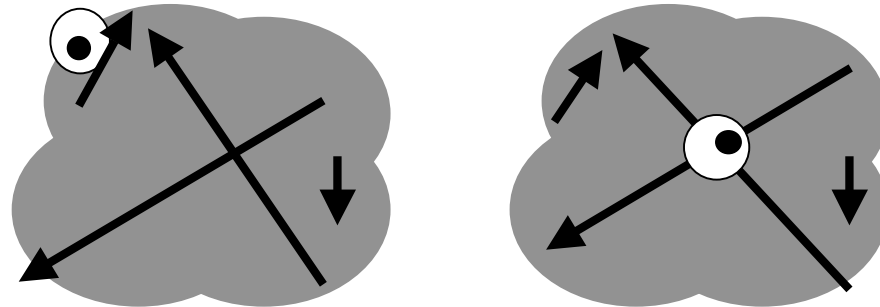


# NetFlows Instrumentation Issues

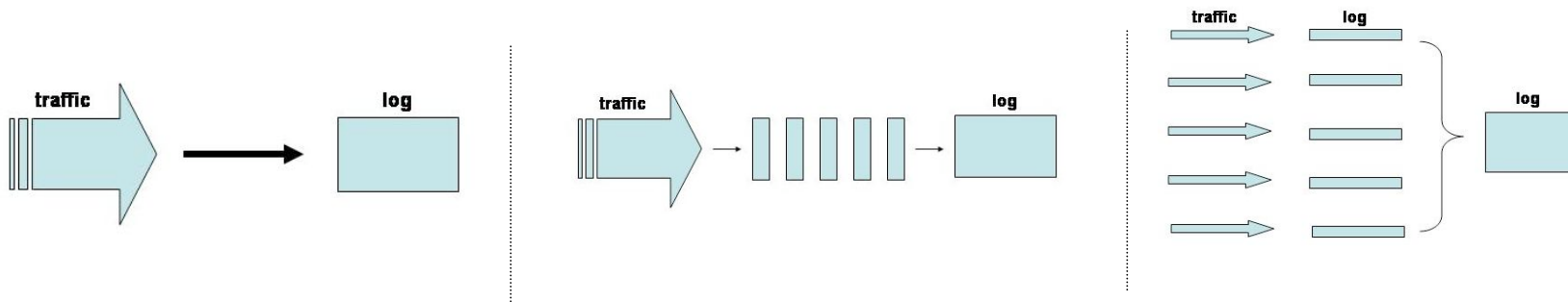
- Streaming Data



- Vantage Point



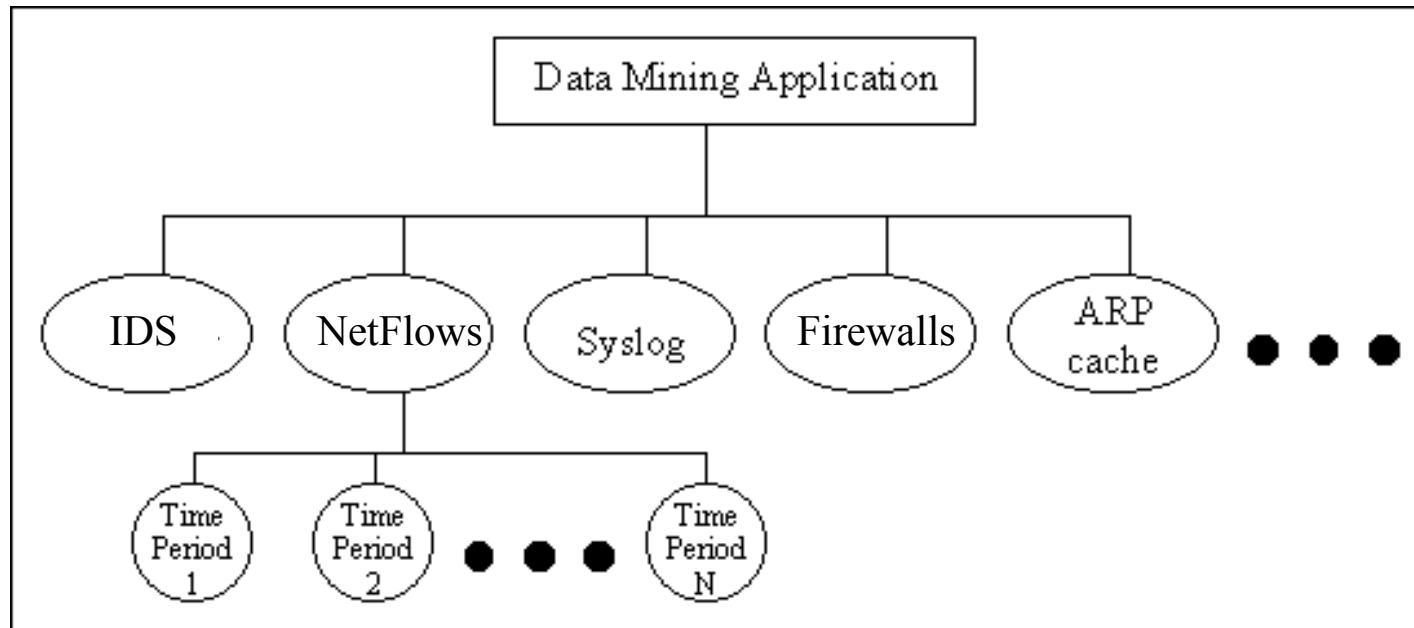
- High Line Rates



# Flavors of NetFlows

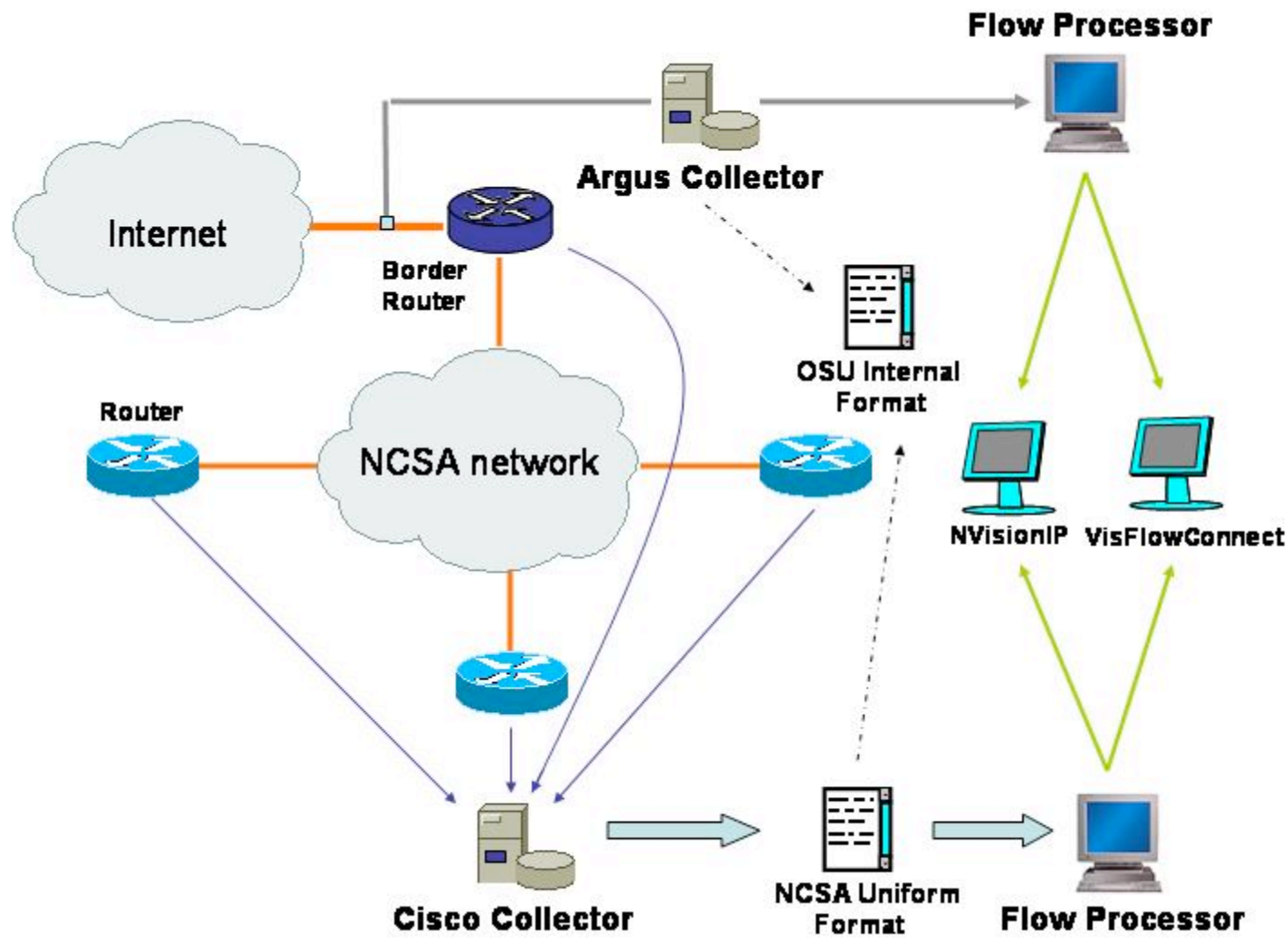
- **Router-Based (Cisco, Juniper, etc.)**
  - Cache timeout
  - Configuration
  - Sampling
- **Argus** <[\*http://www.qosient.com/argus/\*](http://www.qosient.com/argus/)>
  - Open Source
  - Platform Independent
  - Configuration (data field)
- **Home Grown NetFlows**
  - Many, for instance, Tom Daniels, Iowa State University .....

# The Data Management Problem

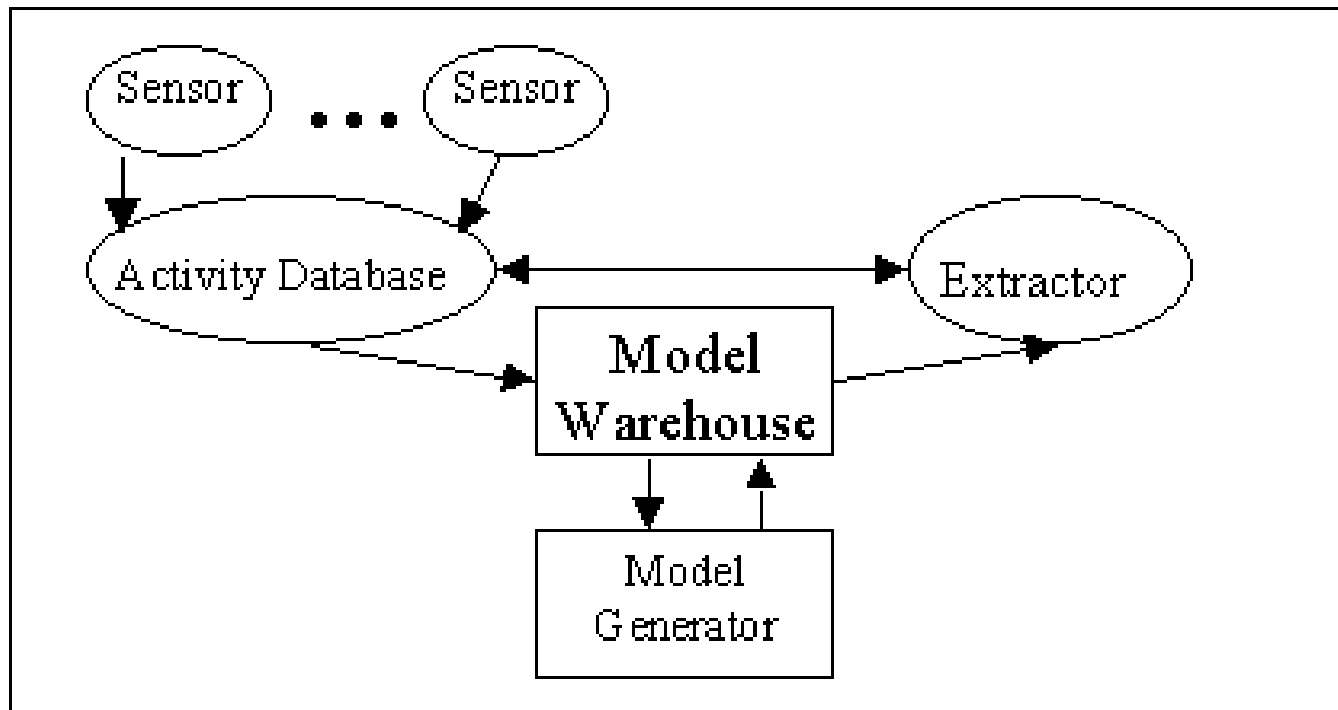


*time dimension*

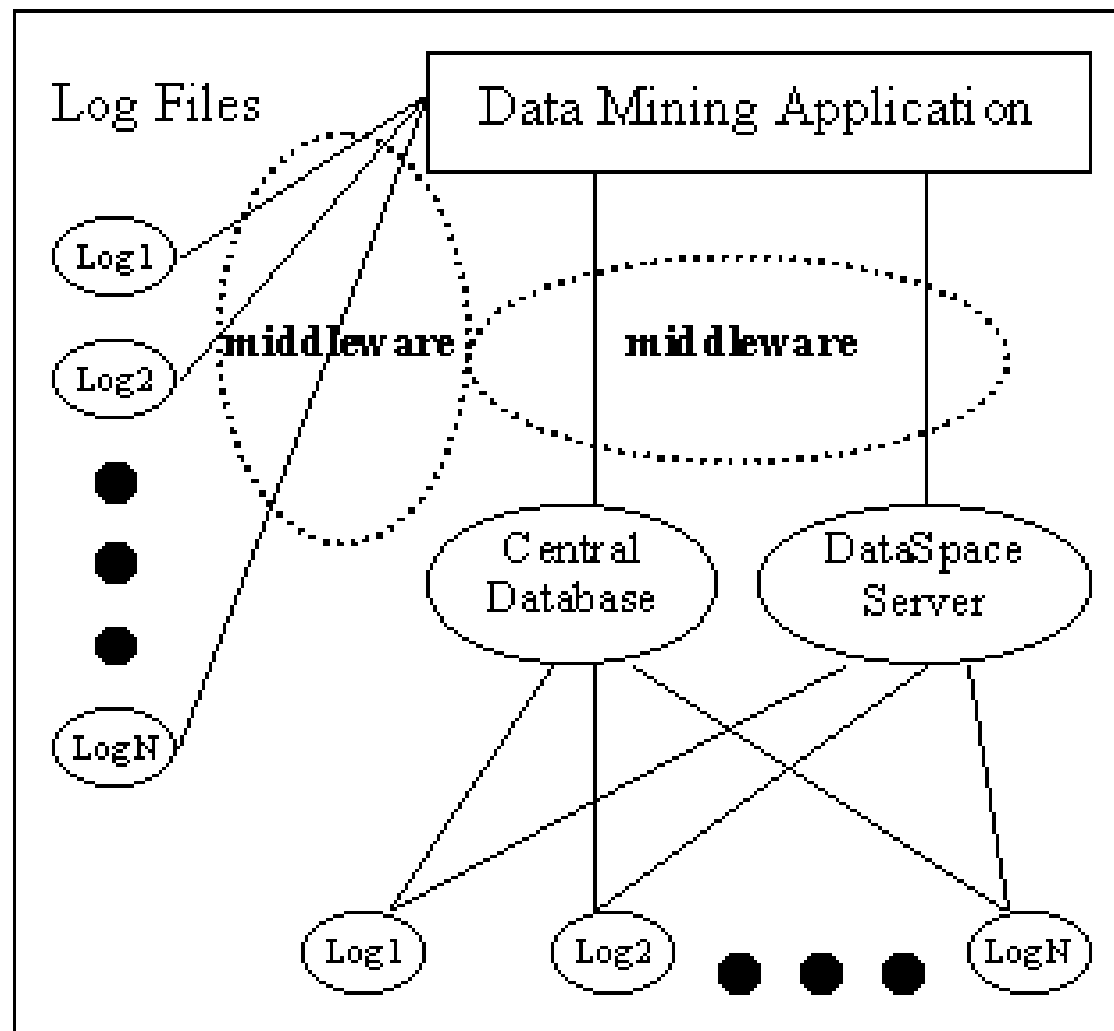
# NCSA's NetFlows Architecture



# (1) Central Database Architecture



## (2) Middleware Architecture





**Carnegie Mellon  
Software Engineering Institute**

**CERT**  
Situational  
Awareness

# Data Sharing: Lessons learned by the CERT/CC and the CERT/NetSA groups

Roman Danyliw <rdd@cert.org>

FloCon 2004: Data Sharing Panel

CERT® Network Situational Awareness Group  
Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

*The CERT Network Situational Awareness Group is part of the Software Engineering Institute. The Software Engineering Institute is sponsored by the U.S. Department of Defense.*





## Background

---

- CERT/CC has a long history of accepting incident reports, artifacts, and vulnerability information
  - Synthesizing this input into public analysis such as advisories and the coordination of patch releases
- CERT/SA has experience in analyzing operational data-sets of other organizations
  - Synthesizing these data-sets to form situational awareness, and new analytical approaches





## Decomposing “Data Sharing”

---

- Data collection
  - Accepting data from outside your organization
- Data dissemination
  - Providing value-add back to data sources or constituency

*An organization only involved in data collection  
is not “data sharing”*



## Concerns in Sharing

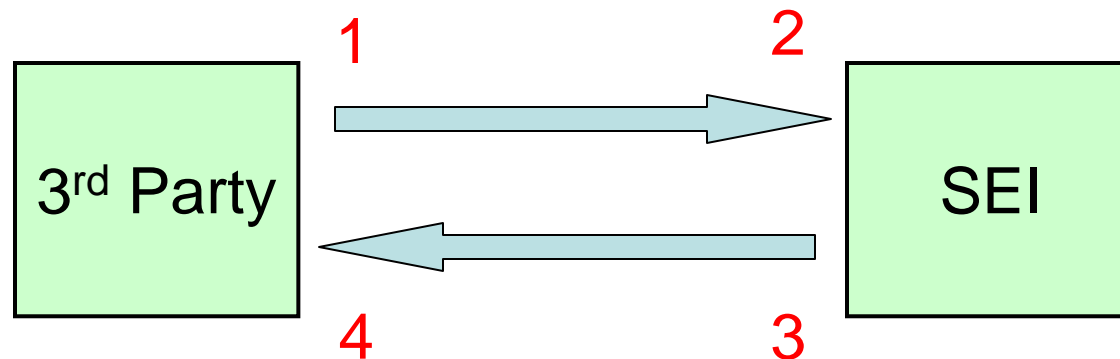
---

- Concerns for the data source
  - Is anything “sensitive” being released?
    - If so, what assurances do I have about my data?
  - Is there sufficient benefit to me in providing this information?
- Concerns for the data recipient
  - Is there any risk in accepting this information?
    - Does the data source know it is a data source?
    - Can others know that this data source is being used?
    - What responsibilities do I have with respect to handling/sharing this information with others?
  - Is there sufficient benefit to collecting this information?



## Steps in the Sharing Process

---





## (1) I am reporting data to CERT

---

- Sharing data is technologically hard and requires human intervention
  - Few tools provide native support for sharing
  - CERT does provide tools to extract, filter, and sanitize information
- What guarantees do I have for my data?
  - Once data is handed over, all guarantees are founded on trust – no practical technological solution
  - Accreditation of processes, technology, and facilities



## (1) I am reporting data to CERT (cont'd)

---

- “My information is sensitive, I want to protect:”
  - Information revealed in packet payloads
    - Contents of email, clear-text authentication
  - Internal topology of the network
    - Size and the purpose of individual hosts
  - Laxness or lapses in security
    - Outbound attacks
    - Usage of certain services (e.g., P2P)
    - Indications of vulnerabilities
- Often raw data is not possible; only share summaries



## (2) CERT is receiving my information

---

- Willingness to share does not always mean utility for the CERT
  - Impossible to mechanically parse free-form text reports
  - Organizational or obscure data formats (i.e., vendor X with tool Y version Z.zzz.z)
- Employ standard data use policies
  - For all automated data sharing, a formal MOU governs the exchange
  - Public, default data disclosure policy for all self-reported data
- Public knowledge of honey-pot addresses is problematic



## (2) CERT is receiving my information

---

- Community specific constraints
  - Academic community
    - Cannot tie data back to students
      - IP address resolved to host names which contained a student's name
  - Federal community
    - Cannot collect Personally Identifiable Information (PII)
      - Only present in the payload
  - Medical community
    - HIPPA prevents PII collection
      - Only present in the payload



## (3) CERT is disseminating information

---

- Does not provide attribution
  - Sometimes obfuscates results to do peer comparison
- Coordinating pre-release information requires a substantial volume of encrypted email
  - Dedicated tool (srmil) to handle encryption/decryption among various standards (e.g., gpg, pgp, s/mime)
- How to control the use of data after it is made available?
  - Contractors and federal government “rights to use” on pre-release information
  - Data leak through a 3<sup>rd</sup> party
  - Reaction of some open-source vs. COTS vendors to a vulnerability





## (3) CERT is disseminating information

---

- Who is the right audience?
  - Traditionally, advisories were for system administrators – now have summaries for management
  - How to reach home users?



## (4) I am receiving CERT information

---

- Optimal format for receiving information:
  - Semantics: push vs. pull
  - Transport protocol: email, web, etc.
  - Machine parsable vs. human readable
- How timely is the information?
  - Incomplete information, but early notification
    - Incremental updates
  - Complete information, but late notification



## Observations in Data Sharing

---

- Datasets based on more sites is not always better – a representative sample is key
  - *Defining representative is hard*
- The community needs to develop and adopt standards formats and protocols to exchange analytical results
  - *Adoption by the vendor community will be required*
- Centralization is not desirable; expertise to analyze data is rarely found in one place – build a community of analysts
  - *The politics of data sharing make this hard*

## ***Sharing Intelligence is our Best Defense:***

### **Incentives That Work versus Disincentives That Can Be Solved**

William Yurcik\*   Adam Slagell   Jun Wang

NCSA Security Research  
National Center for Supercomputing Applications (NCSA)  
University of Illinois at Urbana-Champaign

*Data Sharing Panel  
FloCon 2004*

National Center for Supercomputing Applications



## ***Cyber Security Today Is “a bit” Like the Keystone Cops***



National Center for Supercomputing Applications



## ***Cyber Security Today Is “a bit” Like the Keystone Cops***



*They do  
something  
really bad!*

National Center for Supercomputing Applications



## ***Cyber Security Today Is “a bit” Like the Keystone Cops***



*They do  
something  
really bad!*

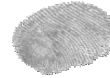
*Then we chase  
them to the  
border.*

National Center for Supercomputing Applications



## Security Information Sharing

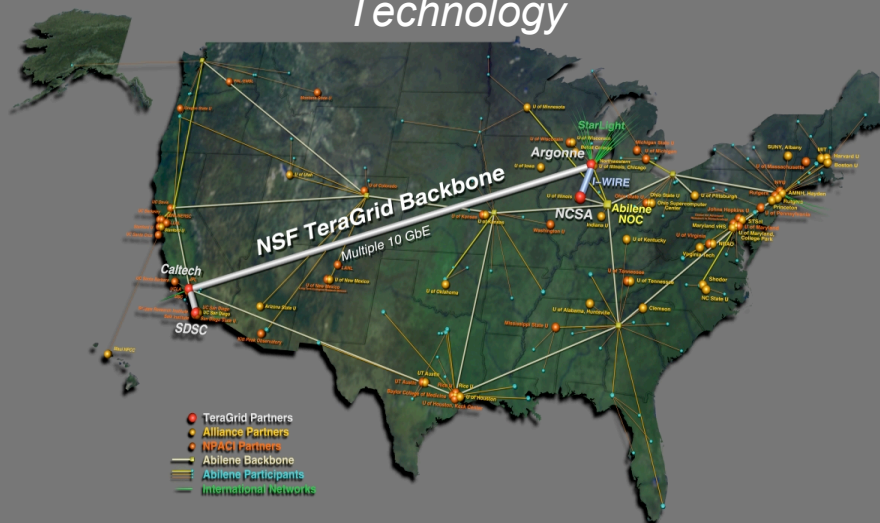
- Need to share information on attacks.
  - Fingerprints and attack profiles
  - Individual events
- Identify individuals
- We cannot continue to stop at the border, we need to cooperate with law enforcement and each other.
  - Security event repository
  - Event correlation across administrative domains
- “unfortunately, this country takes **body bags** and requires **body bags** sometimes to make really tough decisions about money and about governmental arrangements” - Richard Clarke 9/11 Testimony



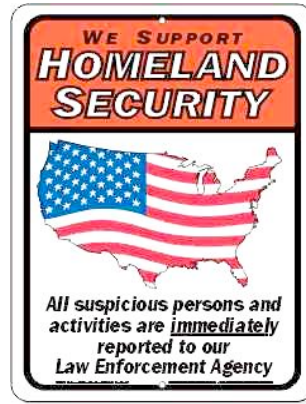
National Center for Supercomputing Applications



## *The World is Rapidly Changing Greater Dependency on Collaborations and Technology*



## Cooperation is Voluntary



*The vast majority of incidents are never reported*

National Center for Supercomputing Applications



## Cooperation is Voluntary Caveat - except in California!

*Only state mandatory disclosure law currently on the books at state level.  
Effective as of July 2003*

National Center for Supercomputing Applications



## Cooperation is Voluntary Caveat - except in California!

*Only state mandatory disclosure law currently on the books at state level.  
Effective as of July 2003*

### California Law has national effects:

*California is home to many of the biggest technology companies in the country.*

*Law applies to all who "conduct business" in the state. Of course many companies route their information through servers housed in California.*

*Potential for litigation in California - many times companies will have no way of knowing whether a person is resident of California or not.*

National Center for Supercomputing Applications



## Computer Emergency Response Teams CERTs

<http://www.first.org/team-info/>



National Center for Supercomputing Applications





## Information Sharing and Analysis (ISACs)

- Gathering, analysis and sharing of information related to actual or unsuccessful attempts at computer security breeches.
- Presidential Decision Directive (PDD)-63
- Fee base membership
- Operational ISACs
  - Electric power
  - Telecommunications
  - Information technology
  - Financial services
  - Water supply
  - Surface transportation
  - Oil & gas
  - Emergency fire services
  - Food
  - Chemicals industry
  - Emergency law enforcement

National Center for Supercomputing Applications



**Question:**

**Can we share?**

National Center for Supercomputing Applications

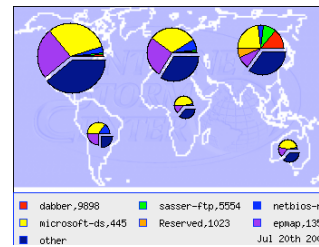
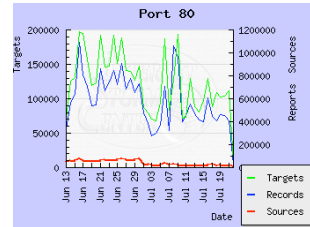
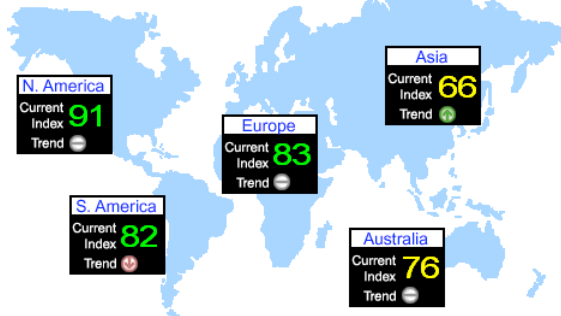




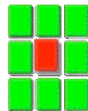
## (1) SANS

**INTERNET  
TRAFFIC  
REPORT**

Last update (MST):  
7/21/2004 20:20  
Global  
Index **85**  
Trend

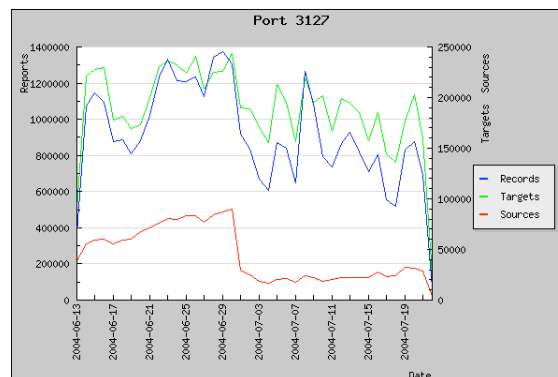


National Center for Supercomputing Applications



## (2) DShield.org

Distributed Intrusion Detection System



Services registered for this port      Vulnerabilities for this port

....

...

National Center for Supercomputing Applications





<<http://www.first.org/>>

### **(3) Forum of Incident Response and Security Teams**



### **(4) CIC-SWG**

**Committee on Institutional Cooperation**

**- IT Security Working Group**

**(Big Ten Universities plus the University of Chicago)**

<<http://www.cic.uiuc.edu/groups/ITSecurityWorkingGroup/>>

National Center for Supercomputing Applications



# **Incentives / Disincentives**

National Center for Supercomputing Applications



## Framing the Data Sharing Issues

- **Both an Internal / External Issue** (within before between)
- **Who should share externally?**
  - at what organizational levels (more/less bureaucracy)
  - flat or hierarchical (scalability)
- **What should be shared?**
  - raw data, processed data, known answers
- **How should it be shared?**
  - phone calls/Emails, reports, automation 😊
- ❌ **Significant time and effort to share**
  - payback? none/long-term ❌ real-time 😊
- ❌ **Does technology exist to share securely**
  - Will information I share come back to bite me?

National Center for Supercomputing Applications



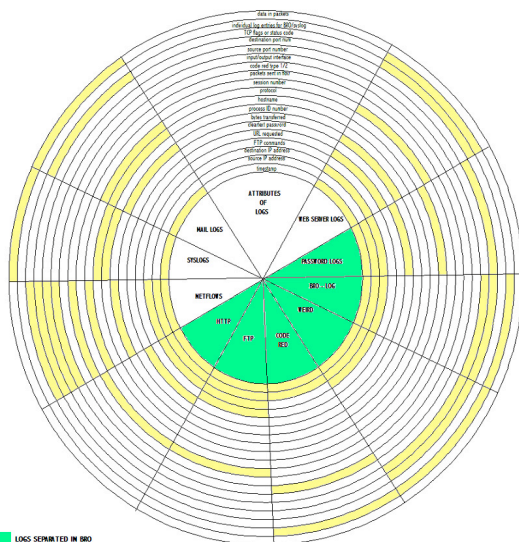
## Commonly Available Logs

- |                                  |                             |
|----------------------------------|-----------------------------|
| 1) <b>NetFlows Logs</b>          | 12) Vulnerability Scan Logs |
| 2) Packet Traces - tcpdump       | 13) Nameserver DNS Cache    |
| 3) Network IDS- BRO, Snort, etc. | 14) SNMP Logs               |
| 4) Host IDS – Tripwire, etc.     | 15) BGP Tables              |
| 5) Syslogs (general)             | 16) Dial-Up Server Logs     |
| 6) Authentication Logs           | 17) ARP Cache               |
| 7) DHCP Server Logs              | 18) Workstation Logs        |
| 8) Firewall logs                 | 19) Process Accounting Logs |
| 9) Mail Server Logs              | 20) Trace Route Logs        |
| 10) Backup Logs                  | 21) “Homegrown” Logs        |
| 11) AntiVirus Logs               | .....                       |

National Center for Supercomputing Applications



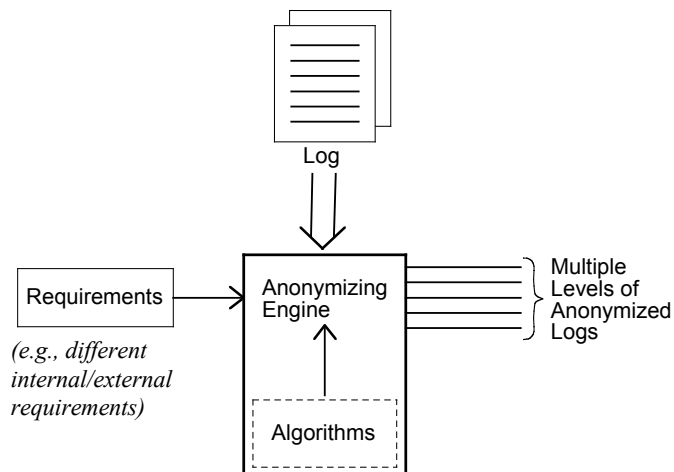
## Attributes Across Logs



National Center for Supercomputing Applications



## Log Anonymization

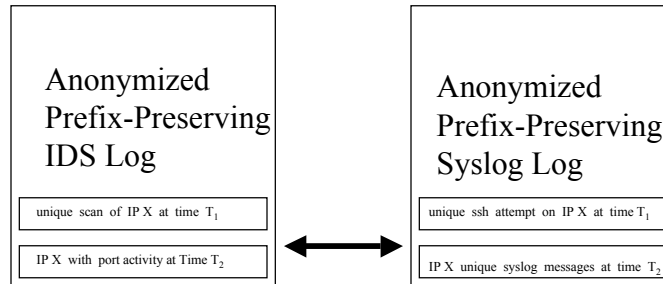


National Center for Supercomputing Applications

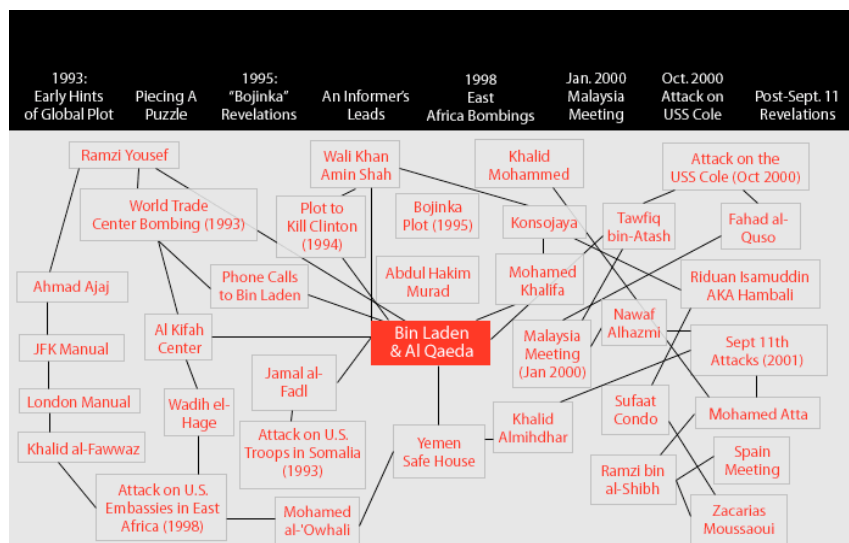


## Known Plain-Text Attacks

## Statistical Inference



National Center for Supercomputing Applications



National Center for Supercomputing Applications



## NCSA SIFT Project

<http://www.ncassr.org/projects/sift/>

VizSEC Workshop Oct 29, 2004  
ACM Computer and Communications  
Security Washington DC  
<http://www.cs.fit.edu/~pkc/vizdmsec04/>

National Center for Supercomputing Applications



## Discussion

- No *one-size-fits-all* solution exists for log sharing
- Solutions depend on the application
  - three major problems
    - 1) huge distributed data volumes
      - visualization is part of the solution here – next workshop
    - 2) security must be considered
      - CIA
      - may require re-design/re-architecture (I hope not!)
    - 3) Incentives
- Operational incentives may be the key
  - We have a counter-intuitive example that actually works:
    - sharing between very selfish sysadmins with very sensitive security information (go figure)
  - “only cooperation will make us less vulnerable”

National Center for Supercomputing Applications





Carnegie Mellon  
Software Engineering Institute

**CERT**  
Situational  
Awareness

# The State of Standardization Efforts to support Data Exchange in the Security Domain

Roman Danyliw <rdd@cert.org>

FloCon 2004: Standards Talk

CERT® Network Situational Awareness Group  
Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

*The CERT Network Situational Awareness Group is part of the Software Engineering Institute. The Software Engineering Institute is sponsored by the U.S. Department of Defense.*







# Overview

---

- Flow and Packet Formats
- Alert and Event Formats
- Context-relevant Formats



# Dimensions in Representation

---

- Usage of representation
  - Transport vs. analysis vs. storage vs. archive
- Volume of data informs representation choice
  - Raw vs. Summaries
    - Choice often dictates a binary vs. text implementation
- Policy Scope
  - Intra-Organizational
    - Little consensus from outsiders necessary
    - Interoperation focus
  - Inter-Organizational
    - Privacy issues more acute (sanitization, filtering)
    - Common semantics are more relevant
    - Efficiency of representation is more significant



## Formats of interest

---

- Flow and Packet Formats
  - IPFIX
  - PSAMP
- Alert and Event Formats
  - IDWG
  - INCH
- Context-relevant Formats
  - Vulnerability Report
  - CRISP



## Flow and Packet Formats (*de facto*)

---

- PCAP (tcpdump)
  - <http://www.tcpdump.org>
- Cisco NetFlow



# IETF IP Flow Information Export (IPFIX) WG

---

<http://www.ietf.org/html.charters/ipfix-charter.html>

- Binary, extensible information model for IP flows exported from a given *observation point* (i.e., router line-card) to a *collector*
  - Based on Cisco Netflow v9
- Designates a mandatory protocol (SCTP) to use in the transport of these flows

(Note: Various text and figures were taken from the IPFIX I-Ds)



## IPFIX Flow Definition

---

- “... [A] set of IP packets passing an observation point ... during a certain time interval. All packets belonging to a particular flow have a set of common properties [named flow keys].”
  - One or more packet header field (e.g. destination IP address), transport header field (e.g. destination port number), or application header field (e.g. RTP header fields)
  - One or more characteristics of the packet itself (e.g. number of MPLS labels)
  - One or more fields derived from packet treatment (e.g. next hop IP address, output interface)



## IPFIX Flow Definition

---

(2)

- A flow is defined by a *flow type* function that considers the various *flow keys*
- Flexible definition provides support for:
  - Filtering
  - Sampling
  - Bi-directional and unidirectional flows



# IPFIX Information Model

---

- Template-based format
  - IPFIX merely specifies the possible
    - data types (e.g., IPv4 address, octet) and the
    - information items (e.g., icmpTypeCode, egressInterface)
  - Information items are unique identifiers registered with IANA or escaped via a vendor code
  - A template is merely an ordered list of pairs:  
<information items (i.e., fieldID), data length>
    - No static format; can be dynamically generated during the export process





# IPFIX Information Model

---

(2)

- Two classes of records
  - Template Records
    - Describe a format
  - Data Records
    - Contain data encoded and formatted according to a Template record
- Two flavors of Data Records; those that encode the:
  - Data stream (e.g., observed flows), and
  - Control Information (e.g., selection criteria)



# IPFIX Information Model

---

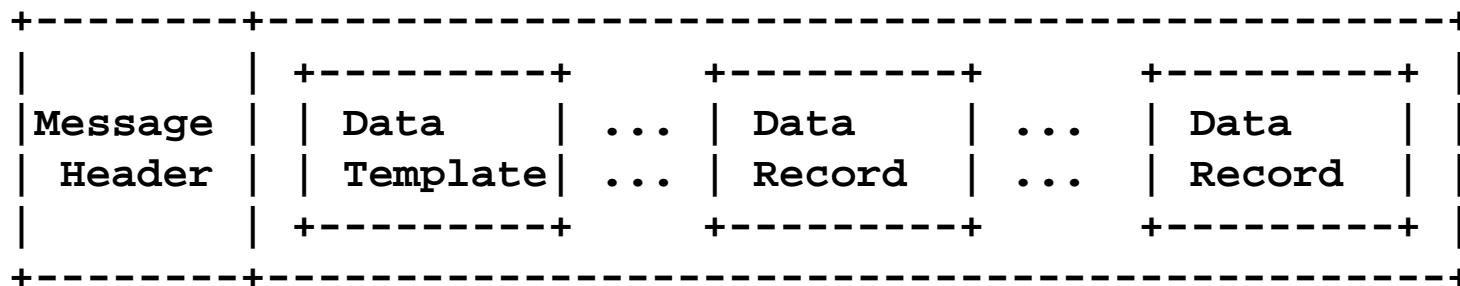
(3)

- 4-basic record types
  - Flow Data Template
    - A description for data record structure
  - Flow Data Record
    - IP flows formatted according to the Flow Data Template
  - Option Template
    - A description of the option record structure
  - Option Record
    - Control information formatted according to the Option Template Record



# IPFIX Messaging

---

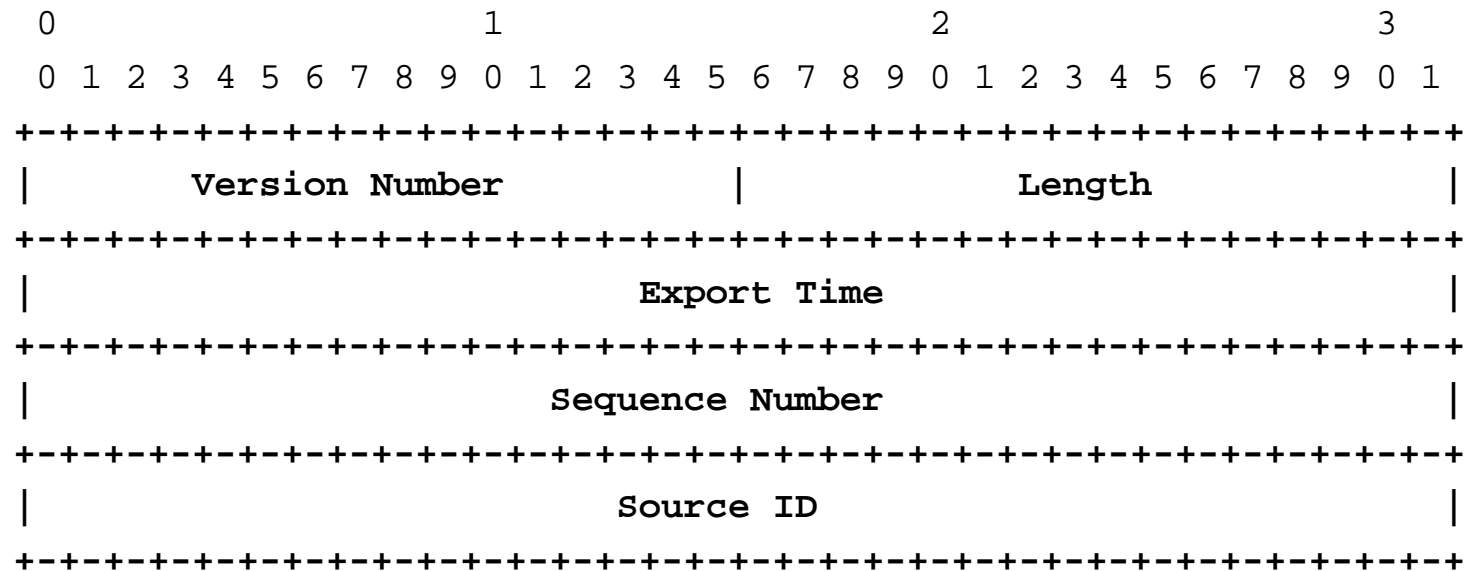


- Template records are sent inline with the data records
  - Frequency dictated by the quality of transport
  - Possible to send no template in an export, and reference a previously sent template in the data record
    - Collector must cache data templates



# IPFIX Message Header

---



- 128-byte preamble sent with each export



# IPFIX Example

---

Src IP addr.	Dst IP addr.	Packet Number	Bytes Number
-----	-----	-----	-----
198.168.1.12	10.5.12.254	5009	5344385
192.168.1.27	10.5.12.23	748	388934

Flow  
Information  
to Export

+-----+			
		+-----+	+-----+
Message	Data	Data	
Header	Template	Records	
	(1 Template)	(2 Flow Data Records)	
		+-----+	+-----+
+-----+			

IPFIX  
Encoding  
Format



# IPFIX Example: Template

(2)

0										1										2										3									
0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1								
FlowSet ID = 0										Length = 24 bytes																													
Template ID 256										Field Count = 4																													
IP_SRC_ADDR = 0x0008										Field Length = 4																													
IP_DST_ADDR = 0x000C										Field Length = 4																													
IN_PKTS = 0x0002										Field Length = 4																													
IN_BYTES = 0x0001										Field Length = 4																													



## IPFIX Example: Data

(3)

```

0          1          2          3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          FlowSet ID = 256          |          Length = 36          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          198.168.1.12          | #1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          10.5.12.254          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          5009          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          5344385          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          192.168.1.27          | #2
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          10.5.12.23          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          748          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|          388934          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```



# IPFIX Transport Protocol: SCTP

---

- Reliable service
  - TCP equivalent
- “Partially reliable” service
  - During un-congested periods, all the records marked for deletion under congestion will be reliably delivered
  - During congested periods, the exporter will drop packets to protect the network





## IPFIX I-Ds

---

- Requirements for IP Flow Information Export
  - draft-ietf-ipfix-reqs-16
- Architecture Model for IP Flow Information Export
  - draft-ietf-ipfix-architecture-03
- Information Model for IP Flow Information Export
  - draft-ietf-ipfix-info-03
- IPFIX Protocol Specifications
  - draft-ietf-ipfix-protocol-03



## IETF Packet Sampling (PSAMP) WG

---

<http://www.ietf.org/html.charters/psamp-charter.html>

- Binary, extensible information model for specifying
  - Selection operations (sampling and filtering) on a packet stream, and
  - Packets yielded by the selection operation
- Designates a mandatory protocol (IPFIX) to use in the transport of these packets



## Relationship between IPFIX and PSAMP

---

- PSAMP extends the IPFIX data model
  - A PSAMP data record is an special instance of an IPFIX flow record with different semantics
    - i.e., a flow record with only a single packet
  - Augments the IPFIX data model to support *Selection Process*
- PSAMP reuses the IPFIX transport protocol



## PSAMP Selection

---

- Sampling
  - “Provisioning of information about a specific characteristic of the parent population at a lower cost than a full census would demand”
- Filtering
  - Deterministic selection of packets based on the
    - packet content
    - treatment of the packet at the observation point, or
    - functions operating on the selection state.
- Possible to create schemes combining of both sampling and filtering selections



## PSAMP Sampling

---

- Systematic Sampling (deterministic function)
  - Count-based (spatial packet position; e.g., packet count)
  - Time-based (temporal packet position; e.g., arrival time)
- Random Sampling
  - n-out-of-N
  - Probabilistic
    - Uniform Probabilistic (same probability for each packet)
    - Non-Uniform Probabilistic (probability depends on input)
    - Flow State Probabilistic
  - Sampling probability depends on flow state



# PSAMP Filtering

---

- Match/Mask
  - Apply bit mask to the header or the first N-bytes
- Hashing
  - Apply a hash function to the header or first N-byte
- Packet Features
  - Properties of the packet header
- Router-state selection
  - Properties of the route or packet treatment



## PSAMP I-Ds

---

- A Framework for Passive Packet Measurement
  - draft-ietf-psamp-framework-05
- Sampling and Filtering Techniques for IP Packet Selection
  - draft-ietf-psamp-sample-tech-04
- Packet Sampling (PSAMP) Protocol Specifications
  - draft-ietf-psamp-protocol-01
- Information Model for Packet Sampling Exports
  - draft-ietf-psamp-info-01



## IETF Intrusion Detection WG (IDWG)

---

<http://www.ietf.org/html.charters/idwg-charter.html>

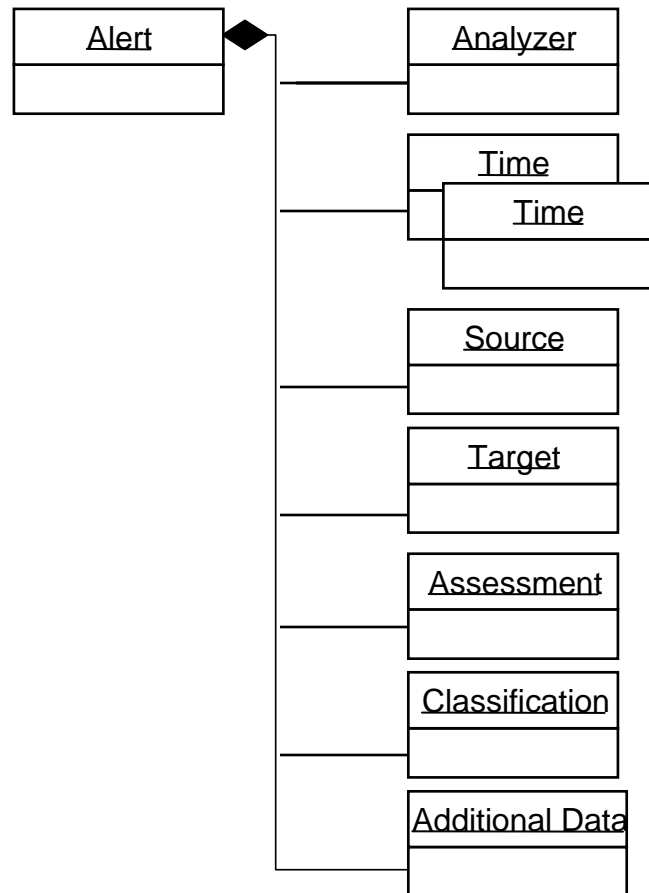
- XML information model for network and host-based Intrusion Detection System alerts
  - Intrusion Detection Message Exchange Format (IDMEF)
- Defines a protocol to exchange these alerts
  - Intrusion Detection Exchange Protocol (IDXP)
  - BEEP-based profile to exchange IDMEF





# IDMEF Data Model

---



- Sensor properties
- Timestamps
- Source/Target characteristics
  - IP address, ports
- Impact assessment
- Event classification
- Extension mechanism



## IDWG I-Ds

---

- Intrusion Detection Message Exchange Requirements
  - draft-ietf-idwg-requirements-10
- The Intrusion Detection Message Exchange Format
  - draft-ietf-idwg-idmef-xml-12
- The Intrusion Detection Exchange Protocol (IDXP)
  - draft-ietf-idwg-beep-idxp-07
- The TUNNEL Profile
  - Rfc3620



## IETF Incident Handling WG (INCH)

---

<http://www.ietf.org/html.charters/inch-charter.html>

- XML information model for exchanging “incident data” among CSIRTs
  - Incident Object Description Exchange Format (IODEF)
- No exchange protocol specified



# INCH IODEF Data Model

---

- Extensible framework to exchange information between CSIRTs
  - Workflow
    - incident identifiers, conveying expectations, data usage restrictions
  - Incident description and conclusions
    - Source/Destination information
    - Contact information
    - References to vulnerabilities, advisories, and artifacts
    - Classification and impact assessments
- Extensions
  - RID: DoS traceback for ISPs
  - (possible) Anti-Spam lists



## INCH I-Ds

---

- Requirements for Format for INcident Report Exchange (FINE)
  - draft-ietf-inch-requirements-03
- The Incident Data Exchange Format Data Model
  - draft-ietf-inch-iodef-02
- The Incident Object Description Exchange Format (IODEF) Implementation Guide
  - draft-ietf-inch-implement-00
- Real-Time Inter-Network Defense
  - draft-ietf-inch-rid-00



## IETF Cross-Registry Information Service Protocol (CRISP) WG

---

<http://www.ietf.org/html.charters/crisp-charter.html>

- XML, extensible information model for global registry information
  - i.e., Whois with structure
- Designates a mandatory protocol (BEEP) for the query/response exchange



## Vulnerability Information (*de facto*)

---

- Mitre CVE
  - <http://cve.mitre.org/>
- Mitre OVAL
  - <http://oval.mitre.org/>
- NIST iCAT
  - <http://icat.nist.gov/icat.cfm>



## Vulnerability (Report) Formats

---

- Common Advisory Interchange Format (CAIF)
  - RUS-CERT
  - <http://cert.uni-stuttgart.de/projects/caif/>
- Advisory and Notification Markup Language (ANML)
  - OpenSec
  - <http://www.opensec.org/anml/>
- Application Vulnerability Description Language (AVDL)
  - OASIS
  - [http://www.oasis-open.org/committees/tc\\_home.php?wg\\_abbrev=avdl](http://www.oasis-open.org/committees/tc_home.php?wg_abbrev=avdl)





## Relevance of the Formats to Flows

---

- IPFIX
  - Storage and transport format for flows
- PSAMP
  - Describe acquisition process of the flows
- IDMEF
  - Describe events created from flows
- IODEF (with/without extensions)
  - Describe flow summaries, baselines, etc.



## Adoption

---

- Packets and Flow Formats
  - IPFIX: implementations exist (e.g., Argus)
  - PSAMP: work in progress
- Alerts and Events Formats
  - IDMEF: adoption only in Snort, Prelude, Arcsight
  - IODEF: adoption by 5-15 CSIRTs in Europe, Asia, and the US
- Context Formats
  - Vulnerability formats: work in progress, some used in closed communities
  - CRISP: work in progress



Carnegie Mellon  
Software Engineering Institute

**CERT**  
Situational  
Awareness

## Wish List

Tom Longstaff

CERT® Network Situational Awareness Group  
Software Engineering Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

*The CERT Network Situational Awareness Group is part of the Software Engineering Institute. The Software Engineering Institute is sponsored by the U.S. Department of Defense.*

© 2004 by Carnegie Mellon University





## Suggestions for the evaluation sheet

---

- Topics for future FloCon
  - Want to help organize the next one?
- Other organizations/individuals that should be involved
- Need for a discussion group (netflo@cert.org)

JTF,

DoE thread in the DoE security conference – add research from DoE labs

U. Of Ill. – cisco wants to be involved. ARGUS developers

CAIDA – program committee format. Not one organization. More operational folks ATT Sprint, etc.

CERT

LE FBI. Nanog – ISP SEC BoF

ACSAC – Dec conference on collecting network data.  
Offered on the Tues of ACSAC

Webcast? – not usually done.

Moving too fast to wait a year.

Tech exchanges with regionals

Good group and open discussions between sessions

Architecture session – everyone gets 5 minutes to present

Ask for a 2 page position paper. Forms a proceedings.



## What do you want from Netflow?

---

- Distribution of flags
- Payload hash
- Start/End packet
- IPV6
- MPLS
- Eval network changes on netflow implementation
- IP packet frags
- Sizing characterization (mean/packets vs packet size)
- Methods of data reduction (sampling, compression, etc)
- ICMP data



## What can be shared?

- What's in the "too hard" category
  - Raw data with specific intrusions against our infrastructure
  - Meta data exposing vulnerabilities on specific machines
- What can be shared *beyond data*
  - Tools
  - Techniques
  - Insights
  - Internet-level activity
  - "normal" indicators

Typical use/sharing policies (library of popular ones)

Obfuscation techniques (without destroying the utility of analysis – may be in the too hard list).

Bad lists/top hit lists/top "n" lists

Spikes (security portal for port/packet volumes)

Queries (to an oracle)

Representations of what you're not seeing (filter rules)

Algorithms we can run on summary data (merging data queries)

How to share others data that shows up in your collection.

Time factors (real time vs historical)

Common modal failures

Attribution information (whois type data) meta data – perhaps all who have the need can make a contribution.

Accurate Geo-location information (techniques, methods, etc – false flag issue, practice, triangulation of multiple sources).

Visualization tools – tool demo session (for workshop?)

CISCO efforts for certs in packets (user to service to server authentication) tracked in Netflow.

Practices (where they can be shared)

Protocols to observe (BGP, etc)

Find more protocols of interest – identify the individuals to the group



## What common tools do we all need?

- Pointer to available tools
- Missing tools
- Prototypes of new analysis techniques
- Visualization tools
- Libraries/queries

Place for people to make comments. Usage and other information.

Domain knowledge to reference.

User interface/easily used flow tools. Faster ramp-up time. Automating low-hanging fruit to make easily tracked traffic automated. Build to more sophisticated environments. "65% of all inbound traffic hitting the router is blaster." Reporting tools. Results that can be understood by sponsors. Goal to change the policy. For DoD, DoE, other... need value that netflow is bringing compared to other technologies. Why to invest?

Share to survive the program. ROI and network situational awareness.

Integration of these tools (framework?)

Issues:

New analysts (mchugh to share tutorial). KB pulled together. Language related to netflow terminology (ala google).

Managers up to speed

Making a general case

Secure portal to get results on your data from others.

Set of specific examples of interesting behavior.

Compare and contrast CISCO netflow with ARGUS and others.

Low level analysts quickly trained on other environments (e.g., PNNL)

KB for analyst to share ideas and results (success results and how that was achieved).

Understand the impact.

Pictures. Marketing strategy for selling netflow results. Translation to english.

People (sharing)?



## What common data do we all need?

---

- Meta-Data
  - Whois historical data
  - Routing data (BGP, traceroute, etc)
- “Normal” samples
  - Traffic patterns and samples from standard servers
  - Normal workstation traffic
  - Acceptable scanning/walking activity
- Internet-wide intrusive activity
  - Flows sampled from popular worms
  - Other flows representing identified behavior

People sharing leads to much greater data trust.

Examples of interesting behavior – scripts and visualizations.

How to continue conversations with legal? Need a practical context to ask questions. Test cases. Come up with an interesting scenario, but be less abstract.